

# Type Theory for Logic and Mathematics

A philosophical and mathematical introduction to type theory

Kevin Davey

July 9, 2023



# Contents

<b>I</b>	<b>Logic</b>	<b>5</b>
<b>1</b>	<b>Introduction to the <math>\lambda</math>-calculus.</b>	<b>7</b>
1.1	Preamble . . . . .	7
1.2	Functions and $\lambda$ -abstraction. . . . .	7
1.3	The untyped $\lambda$ -calculus . . . . .	10
1.4	$\alpha$ -equivalence . . . . .	17
1.5	Substitution Re-examined . . . . .	20
1.6	$\beta$ -reduction and Equivalence . . . . .	24
1.7	$\eta$ -reduction and Extensionality . . . . .	31
1.8	The untyped $\lambda$ -calculus as a model of computation - I . . . . .	37
1.9	The untyped $\lambda$ -calculus as a model of computation - II . . . . .	40
1.10	Typed Computation. . . . .	42
1.11	Church-style Typed $\lambda$ -calculus . . . . .	43
1.12	Curry-style Typed $\lambda$ -calculus . . . . .	50
1.13	The Bound-typed $\lambda$ -calculus. . . . .	60
1.14	Other Resources . . . . .	65
1.15	Appendix . . . . .	65
<b>2</b>	<b>Intuitionistic Propositional Logic and the Curry-Howard Correspondence.</b>	<b>83</b>
2.1	Preamble . . . . .	83
2.2	Intuitionistic Logic: Motivation . . . . .	83
2.3	Natural Deduction for Propositional Intuitionistic Logic . . . . .	88
2.4	Some Abbreviations . . . . .	93
2.5	The Curry-Howard Correspondence: Version I . . . . .	95
2.6	Avoiding Misunderstandings: I . . . . .	101
2.7	More Consequences: Proof Normalization. . . . .	103
2.8	Adding Types: Ordered Pairs . . . . .	110
2.9	Do We Need Product Types? . . . . .	114
2.10	Adding Types: Sums . . . . .	116
2.11	Adding Types: The Empty Type . . . . .	121

2.12 $\eta$ Rules for Product and Sum Types. . . . .	122
2.13 Collecting Things Together . . . . .	125
2.14 The Curry-Howard Correspondence: Version II . . . . .	133
2.15 Proof Normalization Revisited . . . . .	138
2.16 The BHK Interpretation Revisited . . . . .	146
2.17 Other Resources . . . . .	148
2.18 Appendix . . . . .	149

**Part I**

**Logic**



# Chapter 1

## Introduction to the $\lambda$ -calculus.

### 1.1 Preamble.

This chapter is an introduction to the  $\lambda$ -calculus, which is a basic tool needed for much of the rest of this book. We begin by presenting the  $\lambda$ -calculus as a language for talking clearly about functions and their application to arguments, and go on to present the  $\lambda$ -calculus as something more like a programming language. We develop some of the basic mathematical ideas of both the typed and untyped  $\lambda$ -calculus. More sophisticated versions of the  $\lambda$ -calculus are presented in subsequent chapters, but the core ideas of these more sophisticated versions remain largely the same as those of the simple version presented in this chapter.

### 1.2 Functions and $\lambda$ -abstraction.

Almost all of mathematics revolves around the manipulation of *functions*. Mathematicians constantly apply functions to inputs to obtain outputs, and combine simple functions in various ways to obtain more complicated functions. Mathematics contains familiar notations for expressing such things. Given a unary function  $f$  and an input  $i$ , we write  $f(i)$  to indicate the value of the function  $f$  at  $i$ . Given two unary functions  $f$  and  $g$ , we write  $f(g(i))$  to indicate the function in which  $g$  is applied to the input  $i$  first, and then  $f$  is applied to the result. Generalizing, combining, and iterating these notations, mathematicians end up writing down complex expressions such as  $f(g(g(x, y), h(z)))$  to denote functions. This notation is perfectly adequate for most ordinary purposes.

Still, in certain situations this sort of notation can be a bit ambiguous. Consider a binary function  $f(x, y)$ . One can think of this function as some sort of process that on two inputs  $i$  and  $j$  produces an output  $f(i, j)$ . But a different point of view is also possible. When a computer applies a binary function  $f(x, y)$  to inputs  $i$  and  $j$ , it might *first* replace  $x$  with  $i$ , and *then*  $y$  with  $j$ . This ordered replacement of  $x$  with  $i$  and then  $y$  with  $j$  is actually a somewhat subtle process. Replacing  $x$  with  $i$ , we end up with a unary function  $f(i, y)$  of  $y$ . We may then replace  $y$  with  $j$

in this expression to obtain the final value  $f(i, j)$ . What we are doing here is treating  $f(x, y)$  as a *unary* function of  $x$ , that takes an input  $i$  for  $x$ , and outputs the unary function  $f(i, y)$  of  $y$ . We can then of course evaluate this function  $f(i, y)$  for some value  $j$  of  $y$ .

Consider for example the function  $f(x, y) = x + y$  on natural numbers. We may think of this as a function of  $x$  that, when input with the number 3, outputs the function  $3 + y$  of  $y$ . Instead of thinking of addition a *binary* function, we are thus thinking of it as a *unary* function that given any number  $x$  as input, outputs the unary function  $x + y$  of  $y$ . If we so desire, we can in this way redescribe any binary function as a unary function.

Of course, we could equally well think of  $x + y$  as a unary function of  $y$  that takes a natural number  $y$  as input, and outputs the unary function  $x + y$  of  $x$ . On this point of view, when input with the number 3,  $f$  outputs the function  $x + 3$  of  $x$ . We thus have two different ways of treating  $x + y$  as a unary function. They will both give us the same result for any ordinary addition problem such as  $2 + 3$ , but these two different ways of treating  $x + y$  as a unary function nevertheless describe two quite distinct processes.

An expression such as  $x + y$  can thus be viewed in several different ways as a unary function. It can be viewed as a unary function of  $x$  that always outputs a function of  $y$ , or as a unary function of  $y$  that always outputs a function of  $x$ . In ordinary mathematical situations this ambiguity does not really matter, and we can just casually think of  $f$  as a binary function of  $x$  and  $y$ . But in other contexts (for example, in describing the step-by-step operation of a computer program) we sometimes must be more careful.

A notation that is very useful for resolving such ambiguities is that of  $\lambda$ -*abstraction*. To indicate that  $x + y$  is to be viewed as a unary function of  $x$ , we use the notation  $\lambda x(x + y)$ , and to indicate that  $x + y$  is to be viewed as a unary function of  $y$ , we use the notation  $\lambda y(x + y)$ . More generally, given any expression  $E$  and a variable  $z$ ,  $\lambda z(E)$  refers to the expression  $E$  viewed as a unary function of  $z$ . (If  $z$  does not occur free in  $E$ , then this will just be a constant function, as we will see later.) Thus the  $\lambda$  operator ‘abstracts’ from an expression, allowing us to view it as a unary function of a particular, single variable.

**Definition 1.1:  $\lambda$ -abstraction.**

If  $E$  is an expression and  $z$  a variable,  $\lambda z(E)$  is the expression  $E$  viewed as a function of  $z$ .

Using this notation, the expression  $\lambda x(x + y)$  is then the expression  $x + y$  viewed as a (unary) function of  $x$ , which on any input  $i$  returns the function  $i + y$  of  $y$ . By contrast, the expression  $\lambda y(x + y)$  is the expression  $x + y$  viewed as a (unary) function of  $y$ , which on any input  $j$  returns the function  $x + j$  of  $x$ .

Most importantly, we can iterate the  $\lambda$  notation to write expressions such as

$$\lambda x(\lambda y(x + y)). \tag{1}$$

Because  $\lambda y(x + y)$  is the expression  $x + y$  viewed as a function of  $y$ , we can then think of (1) as just

$$\lambda x(x + y \text{ viewed as a function of } y)$$



which we can in turn think of as

$(x + y$  viewed as a function of  $y)$  viewed as a function of  $x$ .

This is a function of  $x$ , which on input  $i$ , returns  $i + y$  viewed as a function of  $y$ . This of course is just the way of thinking about addition described earlier, in which  $x + y$  is viewed as a unary function of  $x$  that returns a unary function of  $y$ . The expression

$$\lambda y(\lambda x(x + y)) \tag{2}$$

is, on the other hand, a way of thinking about addition in which  $x + y$  is viewed as a unary function of  $y$  that returns a unary function of  $x$ . In a sense, (1) and (2) describe the same function implemented in two different ways. Thus, although (1) and (2) are in some sense the same function, our new notation allows us to distinguish a point of view from which they are different. This will be useful in what follows.

To improve readability, we will often drop parentheses and just write  $\lambda x \lambda y(x + y)$  or  $\lambda y \lambda x(x + y)$  instead of (1) and (2) where there is no risk of ambiguity.

Some further terminology will also be useful. Just as in ordinary quantificational logic we talk about free and bound variables, we also say that in the expression  $\lambda y(x + y)$ , the variable  $x$  is *free* and the variable  $y$  *bound*. More specifically, we say that the bound variable  $y$  that occurs in  $\lambda y(x + y)$  has been *bound by the operator*  $\lambda y$ . An expression with no free variables is *closed*, and an expression with free variables is *open*. So for example,  $\lambda x \lambda y(x + y)$  is closed, while  $\lambda x(x + y)$ ,  $\lambda y(x + y)$  and  $x + y$  are open. Although the notions of free and bound variables can be given more rigorous definitions, an intuitive understanding of them will suffice in what follows.

### Exercises for Section 1.2

1. The expression  $\lambda x(x \times y)$  has the form  $\lambda x(E)$ , and is thus a unary function of  $x$ . What is its output when 7 is substituted for  $x$ ? (Here ' $\times$ ' refers to multiplication.)
2. What is the output of the function  $\lambda y(x \times y)$  when given the input 7?
3. Using  $\lambda$ -notation, write a unary function that when given an input  $x$ , outputs  $x \times 2$ .
4. Using  $\lambda$ -notation, write a unary function that when given an input  $x$ , outputs a unary function that when given an input  $y$ , outputs  $x \times y$ .
5. Using  $\lambda$ -notation, write a unary function that when given an input  $x$ , outputs a unary function that when given an input  $y$ , outputs a unary function that when given an input  $z$ , outputs  $x + (y \times z)$ .
6. Are the following terms open or closed? What are their free variables, if they have any?

(i)  $\lambda y(x \times y)$

(iii)  $\lambda x \lambda y \lambda z(x + y)$

(ii)  $\lambda x \lambda y \lambda z(x + y + z)$

(iv)  $\lambda x \lambda y(x + y + z)$

### 1.3 The untyped $\lambda$ -calculus

In this section, we introduce the simplest version of the  $\lambda$ -calculus. The  $\lambda$ -calculus comes in many different versions, and we consider several throughout this book. The first version we will consider is an example of a so-called *untyped  $\lambda$ -calculus*. (What the word ‘untyped’ is doing here will only become clear later on when we introduce the ‘typed’ lambda-calculus with which it can be contrasted.) The first version of the untyped lambda calculus we will consider consists of expressions such as

$$\begin{aligned} &\lambda x(\lambda y(x(x(y)))) \\ &\quad x(\lambda x(y(x))) \\ &\lambda y(x(\lambda x(y(x)))) \end{aligned} \tag{3}$$

These expressions are all known as  $\lambda$ -terms. Assuming we have some list  $x, y, \dots$  of variables, the rules of formation for  $\lambda$ -terms are as follows:

**Definition 1.2: Simple untyped  $\lambda$ -terms.**

- (i) Any variable  $x, y, z, \dots$  is a  $\lambda$ -term.
- (ii) If  $M$  and  $N$  are  $\lambda$ -terms, then so is  $M(N)$ .
- (iii) If  $M$  is a  $\lambda$ -term and  $x$  a variable, then  $\lambda x(M)$  is a  $\lambda$ -term.

You should note how these rules can be combined to generate terms such as those in (3). For example, the variables  $x$  and  $y$  are  $\lambda$ -terms by (i). Applying (ii) gives that  $y(x)$  is a  $\lambda$ -term. Applying (iii) then gives that  $\lambda x(y(x))$  is a  $\lambda$ -term, and applying (ii) again gives that  $x(\lambda x(y(x)))$  is a  $\lambda$ -term.

Note that operations like ‘+’ and ‘ $\times$ ’ that appeared in the previous section are not officially part of our language, and thus do not appear in simple untyped  $\lambda$ -terms. While for explanatory purposes it was useful earlier to write expressions such as  $\lambda x(\lambda y(x + y))$ , such terms are not officially part of the simplest version of the untyped  $\lambda$ -calculus we are considering here.

Those who have seen the  $\lambda$ -calculus before will note that the conventions for parentheses given here are different from some other texts. We will discuss this in a little more detail shortly.

What do our  $\lambda$ -terms refer to? We will think of all  $\lambda$ -terms as *unary functions*. They are unary functions which, when given an input – which will also always be a unary function – output a unary function. This unary function can in turn be applied to another input – which will itself be a unary function – yielding yet another unary function as output.

More specifically, given two  $\lambda$ -terms  $M$  and  $N$ , we think of  $M(N)$  as the result of applying the unary function  $M$  to the unary function  $N$ . Rule (ii) above states that whenever we have two  $\lambda$ -terms  $M$  and  $N$ , we can apply the function  $M$  to  $N$ , obtaining  $M(N)$ . Now  $M(N)$  is itself a  $\lambda$ -term and thus a unary function, and so using rule (ii) again it can be applied to any  $\lambda$ -term  $O$  to obtain  $M(N)(O)$ . This  $\lambda$ -term can likewise be applied to any  $\lambda$ -term  $P$  to obtain  $M(N)(O)(P)$ , and so on.

This sort of behavior might seem odd. One might expect such a process of applying functions to inputs to eventually ‘bottom out’ with some sort of value or result which is *not* itself a function.

But in the simplest version of the  $\lambda$ -calculus that we are considering here, computations never ‘bottom out’ in anything that is not a function – all  $\lambda$ -terms are unary functions, the outputs of which are unary functions, and the outputs of *those* unary functions are further unary functions, and so on and so forth, ad infinitum.

On this point of view, even our variables  $x, y, z, \dots$  will be viewed as unary functions. So if  $x$  and  $y$  are variables, then they are both  $\lambda$ -terms and thus unary functions, and the result of applying the unary function  $x$  to the input  $y$  is just the  $\lambda$ -term  $x(y)$ . We can iterate this sort of process using rules **(i)** and **(ii)** to obtain  $\lambda$ -terms  $x(y), x(x), x(y)(z), x(y(z)), x(x(y)(y))$ , each of which denotes a distinct unary function.

Let us turn our attention to rule **(iii)** and the associated concept of  $\lambda$ -abstraction. As an example, consider a term like  $x(y)$ . Using rule **(iii)**, we can perform two different types of ‘abstraction’ over this term, by viewing it as a function of  $x$ , and writing  $\lambda x(x(y))$ , or by viewing it as a function of  $y$ , and writing  $\lambda y(x(y))$ . As an exercise, let us consider the meaning of each of these expressions in turn.

We begin with the expression  $\lambda x(x(y))$ , in which  $x(y)$  is viewed as a function of  $x$ . This function takes an input  $f$  for  $x$ , and outputs  $f(y)$ . That is to say, the function  $\lambda x(x(y))$  transforms a unary function  $f$  into its value  $f(y)$  at  $y$ . So  $\lambda x(x(y))$  may be described as the function ‘evaluate the function on input  $y$ ’. The process of evaluating a function at  $y$  may thus be redescribed as the application of the  $\lambda$ -term  $\lambda x(x(y))$  to that function.

Consider next the expression  $\lambda y(x(y))$ , in which  $x(y)$  is viewed as a function of  $y$ . This function takes an input  $c$  for  $y$ , and outputs  $x(c)$ . That is to say, the function  $\lambda y(x(y))$  transforms an input  $c$  into  $x(c)$ . So  $\lambda y(x(y))$  may be described as the function ‘apply the function  $x$  to the input’. This is of course just a fancy way of redescribing the function  $x$  itself. The activity of applying a function  $x$  to an input may thus be redescribed as applying the  $\lambda$ -term  $\lambda y(x(y))$  to that input.

Given any  $\lambda$ -term  $M$  and variable  $x$ , rule **(iii)** allows us to form the term  $\lambda x(M)$ , in which  $M$  is viewed as a (unary) function of  $x$ . Typically, this will be most useful when  $M$  contains a free occurrence of the variable  $x$ . However, we do not require  $x$  to appear free in  $M$ , and so rule **(iii)** allows us to construct  $\lambda$ -terms such as  $\lambda x(y(y))$ . It also allows us to abstract over the same variable twice, as in the  $\lambda$ -term  $\lambda x(x(\lambda x(x(y))))$ . We discuss the meaning of these terms shortly.

Because too many parentheses can sometimes impair readability, we will sometimes drop parentheses where there is no risk of ambiguity or the intended reading is obvious from context. So for example, the expression  $xy$  will be an abbreviation for  $x(y)$  (i.e.,  $x$  applied to  $y$ ), and the expression  $\lambda vE$  will be an abbreviation for the function  $\lambda v(E)$ . To minimize the number of parentheses, we will also assume association to the left - so for example, the term  $wxyz$  will denote  $w(x)(y)(z)$ , and *not*  $w(x(y(z)))$ . The difference between these two sorts of expressions is important. The expression  $w(x)(y)(z)$  is the function  $w$ , applied to  $x$ , the output of which is a function which is then applied to  $y$ , the output of which is a function which is then applied to  $z$ . By contrast, the expression  $w(x(y(z)))$  is the function  $w$ , applied to the result of applying the function  $x$  to the result of applying the function  $y$  to the function  $z$ . You should think about the difference between these two terms and the way in which it is expressed in the placement of parentheses.

As mentioned already, we will sometimes omit parentheses between adjacent  $\lambda$ -abstractions,

and so instead of writing  $\lambda x(\lambda y(E))$  we will sometimes just write  $\lambda x\lambda y(E)$ . Combining these conventions, instead of writing

$$\lambda x(\lambda y(x(x)(y)))$$

we will generally just write

$$\lambda x\lambda y(xxy)$$

For further compactness, we will sometimes express two (or more) consecutive  $\lambda$  abstractions  $\lambda x\lambda y(\dots)$  as just  $\lambda xy(\dots)$ . Thus, our term can also be written

$$\lambda xy(xxy).$$

If you read other texts on the  $\lambda$ -calculus, you will quickly discover that different textbooks use different conventions for parentheses. Older books often also introduce a complex set of conventions involving dots in order to eliminate the need for parentheses altogether. Such conventions are perhaps useful for designing lean, efficient programs for manipulating  $\lambda$ -terms, but as that is not our goal we will let ourselves be guided instead by a desire to produce the most readable expressions possible. In this regard, note that our notational conventions overlap somewhat with ordinary mathematical practice. In particular, when a  $\lambda$ -term  $f$  is applied to an object  $c$  (which is of course another  $\lambda$ -term), the result in our notation is  $f(c)$ , which is a very familiar way of expressing the evaluation of a function.

At first, the  $\lambda$ -calculus might seem bizarre, and one might wonder what sorts of mathematical objects  $\lambda$ -terms really are. Some practice is required to get the hang of how the  $\lambda$ -calculus works – this will be provided in the next sections and the exercises. At the beginning, it is useful just to think of the untyped  $\lambda$ -calculus as a type of highly abstract formalism, and worry later about what sorts of mathematical objects  $\lambda$ -terms denote.

In order to talk a little more abstractly about the  $\lambda$ -calculus, it will help to introduce yet further notation. A  $\lambda$ -term like  $\lambda x(xxyx)$  is  $xyx$  viewed as a unary function of  $x$ , and so one would expect that applying  $\lambda x(xxyx)$  to a term  $t$  should produce  $ttyt$ . More generally, for any  $\lambda$ -term  $c$ , the function  $\lambda x(M)$  applied to input  $c$  should have the value  $M'$ , where  $M'$  is obtained by substituting  $c$  for every free occurrence of  $x$  in  $M$ . We will use the notation  $M[c/x]$  to indicate the expression in which every free occurrence of  $x$  in  $M$  is replaced by  $c$ . (Note: a minority of texts use the notation  $M[x/c]$  for this.) So for example,  $xyx[t/x]$  is  $ttyt$ , and  $xyx[t/y]$  is  $xytx$ . The function  $\lambda x(M)$  applied to input  $c$  will then have the value  $M[c/x]$ . For some more examples, the function  $\lambda x(y(x))$  on input  $c$  is  $y(x)[c/x]$  which is  $y(c)$ , and the function  $\lambda x(x(x))$  on input  $c$  is  $x(x)[c/x]$  which is  $c(c)$ . In summary:

**Definition 1.3: Substitution notation.**

The function  $\lambda x(M)$  applied to  $c$  has the value  $M[c/x]$ , where  $M[c/x]$  is the expression in which every free occurrence of  $x$  in  $M$  is replaced by  $c$ .

There are actually some subtleties involved in this process of substitution, but these will be dealt with in the next sections.

It is important to emphasize that the substitution notation  $M[c/x]$  is just a way for us to indirectly describe  $\lambda$ -terms built up using Definition 1.2. We are *not* here adding a clause to Definition 1.2 that says that if  $M$  and  $c$  are  $\lambda$ -terms and  $x$  is a variable, then  $M[c/x]$  is a  $\lambda$ -term. The notation  $M[c/x]$  is simply notation we use in talking about  $\lambda$  terms, but it is not part of the official syntactic structure of  $\lambda$ -terms themselves.

Manipulating some simple examples of  $\lambda$ -terms can help give a feel for their meaning. Let us begin with a very simple but central example; the term

$$\lambda x(x).$$

On any input  $c$ , this evaluates to  $x[c/x]$ , which is of course  $c$ . Thus  $\lambda x(x)$  is a function whose output is always identical to its input, and so is just the identity function. Note however that this is a very abstract and powerful sort of identity function, which returns itself on literally *any* input. In particular, consider the term:

$$\lambda x(x)(\lambda x(x)). \tag{4}$$

This is the identity function applied to the identity function, which using our substitution rule is easily seen to simplify to the identity function  $\lambda x(x)$  itself. In more detail, because (4) is the function  $\lambda x(x)$  applied to  $\lambda x(x)$ , it has the value

$$x[\lambda x(x)/x]$$

i.e., the expression  $x$  with every free occurrence of  $x$  replaced with  $\lambda x(x)$ . This is just  $\lambda x(x)$ .

Of course, the idea of applying the identity function to itself is problematic from a set theoretic point of view, because insofar as we view a function as a set of ordered pairs, no function can take itself as an input. There is a sense then in which the functions we can define in the untyped  $\lambda$ -calculus are more general than anything that can be defined from a set-theoretic point of view.

For another example, consider the closed term

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

If one includes all parentheses, this term is  $\lambda x(x(x))$ . This term describes a function that takes  $x$  to  $x(x)$  – that is, this term describes a function that takes an expression into that expression composed with itself. To see this more carefully, note that we have:

$$\lambda x(x(x))(f) = x(x)[f/x]$$

$$= f(f)$$

Note however that this is quite different from composing a function with itself in the traditional sense. Given a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  such as  $f : x \rightarrow \sin(x)$ , we can of course consider the function  $g : x \rightarrow \sin(\sin(x))$ , which is, in a sense,  $f$  composed with itself. However, this sort of self-application is *not* what we mean when we write a  $\lambda$ -term like  $f(f)$ . When we evaluate  $\sin(\sin(x))$  in traditional mathematics, we apply the  $\sin$  function to some *number*  $\sin(x)$ , where we are thinking of  $x$  as having been given. We are *not* applying the  $\sin$  function to the  $\sin$  function itself - that would be meaningless, as the  $\sin$  function acts only on elements of  $\mathbb{R}$ , and not on functions. In contrast, when we write  $f(f)$  in the  $\lambda$ -calculus, we are thinking of  $f$  as a function, and are applying this function to the function  $f$  itself. This is a very different (and less familiar) type of self-application. This must be kept in mind to avoid confusion.

Of course, for  $\lambda$ -terms  $f$  and  $x$ , the  $\lambda$ -term  $f(f(x))$  does denote the function  $f$  applied to the output of the function  $f$  on input  $x$ . In this way, we can represent the more familiar type of function composition and iteration in the  $\lambda$ -calculus. We can write the  $\lambda$ -term

$$\lambda x(f(f(x))) \tag{5}$$

to denote the function that takes an input  $c$  into  $f(f(c))$ . To see in more detail that on input  $c$  the expression (5) gives output  $f(f(c))$ , note that we have

$$\begin{aligned} \lambda x(f(f(x)))(c) &= f(f(x))[c/x] \\ &= f(f(c)) \end{aligned}$$

We can also write

$$\lambda y x(y(y(x))) \tag{6}$$

to denote the function which, when provided with the two inputs  $f$  and  $c$  (in that order) produces the output  $f(f(c))$ . To see in more detail that on inputs  $f$  and  $c$  the expression (6) gives output  $f(f(c))$ , note that we have

$$\begin{aligned} \lambda y x(y(y(x)))(f) &= \lambda x(y(y(x)))[f/y] \\ &= \lambda x(f(f(x))) \end{aligned}$$

That on input  $c$  the function  $\lambda x(f(f(x)))$  outputs  $f(f(c))$  has already been shown. Thus, the closed  $\lambda$ -term (6) may be viewed as an abstract description of the process of applying a function to an argument twice. This more familiar type of iteration of a function  $f$  expressed in both (5) and (6) is quite different from the kind of iteration given by  $f(f)$ .

What the example of the unusual  $\lambda$ -term  $f(f)$  brings out nicely is the fact that  $\lambda$ -terms are expressions that may be legally applied to *any* other expression. In this sense, they are quite unlike functions such as  $\sin$  that are only defined on a limited domain. We will later put some restrictions on the  $\lambda$ -calculus to incorporate functions with limited domains, but until we do so,

we must accept that  $\lambda$ -terms are quite unlike the ordinary mathematical functions with which we are familiar.

For a final example, consider the  $\lambda$ -term

$$\lambda x(y(y)). \tag{7}$$

The free variable  $x$  does not appear inside  $y(y)$ , and so this term is a little unusual. If we input some  $c$  into this function, the output is  $y(y)[c/x]$ . This is  $y(y)$  with every free occurrence of  $x$  replaced by  $c$ . Because there are no free occurrences of  $x$  in  $y(y)$ , this is just  $y(y)$ . So the output  $y(y)$  of this function is the same for every input  $x$ . The function (7) is therefore simply the constant function whose output is always  $y(y)$ .

Before moving on, one further definition will be useful. We often want to talk about the *subterms* of a given  $\lambda$ -term. For example, the subterms of

$$x(\lambda x(x(y)))$$

are  $x$ ,  $y$ ,  $x(y)$ ,  $\lambda x(x(y))$  and the term  $x(\lambda x(x(y)))$  itself. Although an intuitive understanding of this notion will generally suffice, it is nevertheless useful to define it formally. To this end, we have the following definition

**Definition 1.4: Subterms.**

The notion of a *subterm* of a  $\lambda$ -term is defined inductively as follows

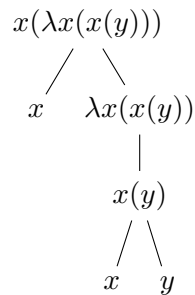
- (i) If  $x$  is a variable, the only subterm of  $x$  is  $x$ .
- (ii) If  $M$  and  $N$  are  $\lambda$ -terms, the subterms of  $M(N)$  are  $M(N)$ , the subterms of  $M$ , and the subterms of  $N$ .
- (iii) If  $M$  is a  $\lambda$ -term and  $x$  a variable, the subterms of  $\lambda x(M)$  are  $\lambda x(M)$  and the subterms of  $M$ .

In the  $\lambda$ -term  $x(y(z))$ , note that the subterm  $y$  is properly contained in the subterm  $y(z)$ . The subterm  $x$  is however disjoint from the subterm  $y(z)$ . In fact, these are in a sense the only possibilities. More specifically, it is impossible for two subterms to only partially overlap. Formally, we have the following:

**Lemma 1.5: Subterm Lemma.**

If  $t_1$  and  $t_2$  are distinct subterms of a  $\lambda$ -term  $s$ , then either  $t_1$  is properly contained in  $t_2$ ,  $t_2$  is properly contained in  $t_1$ , or  $t_1$  and  $t_2$  are disjoint.

The proof is left to the exercises. As a consequence of this lemma, all the subterms of a given  $\lambda$ -term form a treelike structure under the relation of proper containment. For example, the subterms of  $x(\lambda x(x(y)))$  may be represented by the following tree



where a subterm is located beneath another subterm in this tree iff it is properly contained in it. Note that there are two distinct subterms  $x$  of  $x(\lambda x(x(y)))$ . In such a situation, to be fully clear we talk about distinct *occurrences* of the subterm  $x$ . Where there is no risk of confusion we simply talk about subterms simpliciter.

**Exercises for Section 1.3**

1. What is the value of the closed  $\lambda$ -term

$$\lambda x(\lambda y(x(x(x(y))))),$$

when applied to  $x = f$  and  $y = g$ ? Justify your answer in the way shown above, using the  $M[c/x]$  style of notation.

2. What closed  $\lambda$ -term takes expressions  $f$ ,  $g$  and  $h$  (in that order) into the expression

$$(f(g))(f(h))?$$

3. What function does the  $\lambda$ -term  $\lambda x(x)(\lambda y(y))$  describe?

4. (a) Do the closed  $\lambda$ -terms

$$\lambda x(\lambda y(x(y)))$$

and

$$\lambda x(\lambda y(y(x)))$$



describe the same function? (Hint: apply both of these functions to an arbitrary  $f$ , and then apply the result to an arbitrary  $g$ . Are the results the same?)

(b) Do the closed  $\lambda$ -terms

$$\lambda x(\lambda y(x(y)))$$

and

$$\lambda y(\lambda x(y(x)))$$

describe the same function?

5. Are the following true or false? If false, give a counterexample. If true, simply state so.

(a) For any  $\lambda$ -term  $M$ ,  $\lambda x(M)$  is a closed  $\lambda$ -term.

(b) If  $\lambda x(M)$  is a closed  $\lambda$ -term and  $N$  an arbitrary  $\lambda$ -term, then  $\lambda x(M)(N)$  is a closed  $\lambda$ -term.

(c) If  $M$  and  $N$  are closed  $\lambda$ -terms, then  $M(N)$  is a closed  $\lambda$ -term.

6. Prove Lemma 1.5. (Hint: use induction on the construction of  $s$ .)

## 1.4 $\alpha$ -equivalence

We argued earlier that  $\lambda x(x)$  represents the identity function. But there is nothing special about  $x$  here. If  $\lambda x(x)$  represents the identity function, then surely  $\lambda y(y)$  or  $\lambda z(z)$  also represents the identity function. Renaming bound variables should not change a  $\lambda$ -term in any meaningful way. In this and the next section we make this idea more general and precise. In the course of doing so some subtleties about the nature of substitution will arise that require careful attention.

The idea we would like to capture is that when a  $\lambda$ -term  $Y$  can be obtained from another  $\lambda$ -term  $X$  by renaming bound variables, there is some sense in which the  $\lambda$ -terms  $X$  and  $Y$  are the same. In this case, we will say that these terms are  $\alpha$ -equivalent, and write

$$X =_{\alpha} Y.$$

This of course is not yet a fully precise or adequate definition of  $\alpha$ -equivalence. By considering various problematic and unproblematic examples we will move towards a more rigorous definition of  $\alpha$ -equivalence that we will write down shortly.

As a straightforward example, we will have

$$\lambda x(x) =_{\alpha} \lambda y(y). \tag{8}$$

Another simple example is

$$\lambda x(zxw) =_{\alpha} \lambda y(zyw) \tag{9}$$

and a further example is

$$\lambda x(z(x(\lambda y(y)))) =_{\alpha} \lambda u(z(u(\lambda v(v)))) \tag{10}$$

where in (10), two bound variables have been renamed.

We must however be a little careful. Consider the  $\lambda$ -term

$$\lambda x(y) \tag{11}$$

If we replace the bound variable  $x$  with  $y$ , we get the term

$$\lambda y(y) \tag{12}$$

The term (11) is the constant function which always returns the value  $y$ , and (12) is the identity function. We do *not* want to say that the terms (11) and (12) are the same, and thus do *not* want to say that these terms are  $\alpha$ -equivalent.

The problem here is of course that starting with the term (11), we are trying to replace a bound variable  $x$  with a variable  $y$  that already appears free in the original term (11). Nothing like this happens in examples (8), (9) or (10), which are entirely unproblematic.

With this in mind, we define  $\alpha$ -equivalence more carefully as follows. Given a term  $\lambda u(E)$  and a variable  $v$  which does not appear free in  $E$ , we say that the terms  $\lambda u(E)$  and  $\lambda v(E[v/u])$  are  $\alpha$ -equivalent. So for example, we have

$$\lambda x(x) =_{\alpha} \lambda y(x[y/x]) = \lambda y(y)$$

which gives us (8), and

$$\lambda x(zxw) =_{\alpha} \lambda y(zxw[y/x]) = \lambda y(zyw)$$

which gives us (9).

Let us make two further stipulations. First, we stipulate that the relation  $=_{\alpha}$  is transitive. That is to say, we stipulate that for any  $\lambda$ -terms  $x, y$  and  $z$ , if  $x =_{\alpha} y$  and  $y =_{\alpha} z$ , then  $x =_{\alpha} z$ .

Second, we stipulate that replacing a subexpression of a term with an  $\alpha$ -equivalent subexpression gives an  $\alpha$ -equivalent term. More precisely, if a term  $t$  contains a term  $s$  as a subexpression and  $s =_{\alpha} s'$ , then if  $t'$  is the result of replacing the subexpression  $s$  of  $t$  with  $s'$ , then  $t =_{\alpha} t'$ . So for example, because  $\lambda x(x) =_{\alpha} \lambda y(y)$ , we have

$$z(\lambda x(x)) =_{\alpha} z(\lambda y(y))$$

(where we have let  $t = z(\lambda x(x))$ ,  $s = \lambda x(x)$ ,  $s' = \lambda y(y)$  and  $t' = z(\lambda y(y))$ ), and

$$\lambda v(z(\lambda x(x))(x(y))) =_{\alpha} \lambda v(z(\lambda y(y))(x(y)))$$

(where we have let  $t = \lambda v(z(\lambda x(x))(x(y)))$ ,  $s = \lambda x(x)$ ,  $s' = \lambda y(y)$  and  $t' = \lambda v(z(\lambda y(y))(x(y)))$ ).

Armed with all this, we then have for example

$$\lambda x(z(x(\lambda y(y)))) =_{\alpha} \lambda x(z(x(\lambda v(v))))$$

and

$$\lambda x(z(x(\lambda v(v)))) =_{\alpha} \lambda u(z(u(\lambda v(v)))).$$

From transitivity we then have (10).

We thus define  $\alpha$ -equivalence rigorously as follows

**Definition 1.6:  $\alpha$ -equivalence.**

Two  $\lambda$ -terms  $t$  and  $t'$  are  $\alpha$ -equivalent iff there is a sequence of  $\lambda$ -terms  $s_0, \dots, s_n$  such that

- (i)  $t = s_0$  and  $t' = s_n$ , and
- (ii) each  $s_{i+1}$  may be obtained from  $s_i$  by replacing some subexpression  $\lambda v(E)$  of  $s_i$  with the expression  $\lambda w(E[w/v])$ , where  $w$  does not appear free in  $E$ .

Note that we regard every term as a subexpression of itself, so that for any  $\lambda$ -term  $E$ ,  $E$  is a subexpression of  $E$ . Also note that the transitivity of  $=_\alpha$  is not built into this definition, but is an easy proven consequence of it. In fact, it can easily be shown that  $=_\alpha$  is an equivalence relation – that is, it can easily be shown that for any  $\lambda$ -terms  $x, y$  and  $z$ , (i)  $x =_\alpha x$ , (ii) if  $x =_\alpha y$  then  $y =_\alpha x$ , and (iii) if  $x =_\alpha y$  and  $y =_\alpha z$ , then  $x =_\alpha z$ . The argument for this is left to the exercises.

Using this definition we can show (10) by setting

$$\begin{aligned} s_0 &= \lambda x(z(x(\lambda y(y)))) \\ s_1 &= \lambda x(z(x(\lambda v(v)))) \\ s_2 &= \lambda u(z(u(\lambda v(v)))) \end{aligned}$$

It is conventional to regard  $\alpha$ -equivalent terms as notational variants of each other. Thus, it is conventional to think of  $\lambda x(x)$  and  $\lambda y(y)$  as the same  $\lambda$ -term, just written in two different ways. When we talk about  $\lambda$ -terms, we will therefore technically be talking about an equivalence class of terms or expressions, each of which is just a notational variant of the other in which bound variables are renamed. So for example, when we talk about the identity function, we may well write down  $\lambda x(x)$  to notationally represent it, but what we really mean is the *class* of mutually  $\alpha$ -equivalent  $\lambda$ -terms

$$\{\lambda x(x), \lambda y(y), \lambda z(z), \dots\}$$

Because talking about classes of  $\alpha$ -equivalent  $\lambda$ -terms is cumbersome, we will follow tradition by abusing notation, and will generally state our results in terms of individual  $\lambda$ -terms. This cuts down on what would otherwise be an excess of notation. The reader concerned with such details should easily be able to distinguish when we are talking about individual  $\lambda$ -terms and when we are really talking about equivalence classes of  $\lambda$ -terms under the relation of  $\alpha$ -equivalence.

Because we regard  $\lambda x(x)$  and  $\lambda y(y)$  as the same term, we will also regard things out of which they are built – e.g.,  $z(\lambda x(x))$  and  $z(\lambda y(y))$  – as the same term. Thus we will have

$$z(\lambda x(x)) =_\alpha z(\lambda y(y))$$

in accordance with Definition 1.6.

**Exercises for Section 1.4**

1. Are the following true or false?

(a)  $x(x) =_{\alpha} y(y)$

(b)  $x(\lambda x(x)) =_{\alpha} x(\lambda y(y))$

(c)  $x(\lambda x(x)) =_{\alpha} y(\lambda y(y))$

2. Which of the following six expressions are  $\alpha$ -equivalent?

$$\lambda x(x\lambda x(x)) \quad \lambda x(x\lambda y(y)) \quad \lambda x(y\lambda y(y))$$

$$\lambda z(x\lambda z(x)) \quad \lambda z(x\lambda z(z)) \quad \lambda z(x\lambda x(z))$$

Where two or more expressions are  $\alpha$ -equivalent, justify your answer using Definition 1.6.

3. Show that  $=_{\alpha}$  is an equivalence relation. That is, show that for any  $\lambda$ -terms  $x, y$  and  $z$

(i)  $x =_{\alpha} x$

(ii) if  $x =_{\alpha} y$  then  $y =_{\alpha} x$ , and

(iii) if  $x =_{\alpha} y$  and  $y =_{\alpha} z$ , then  $x =_{\alpha} z$ .

4. Show that if  $E$  is a  $\lambda$ -term and  $v$  is a variable that does not appear free in  $E$ , then there is a  $\lambda$ -term  $E'$  such that  $E =_{\alpha} E'$  and  $v$  does not appear at all in  $E$ . (Hint: use induction on the construction of  $E$ .)

5. Consider the following alternative to Definition 1.6

Two  $\lambda$ -terms  $t$  and  $t'$  are  $\alpha$ -equivalent iff there is a sequence of  $\lambda$ -terms  $s_0, \dots, s_n$  such that

(i)  $t = s_0$  and  $t' = s_n$ , and

(ii) each  $s_{i+1}$  may be obtained from  $s_i$  by replacing some subexpression  $\lambda v(E)$  of  $s_i$  with the expression  $\lambda w(E[w/v])$ , where  $w$  does not appear at all (either free or bound) in  $E$ .

Are there terms that are  $\alpha$ -equivalent according to Definition 1.6, but not this definition, or vice versa? Either give an argument that these definitions agree in all cases, or give an example of a  $\lambda$ -term on which they differ.

## 1.5 Substitution Re-examined

One might have worries about the definition of  $\alpha$ -equivalence given in the previous section. In particular, it looks like this definition entails the following

$$\lambda x(\lambda y(x)) =_{\alpha} \lambda y(\lambda y(x)[y/x]) = \lambda y(\lambda y(y)) \quad ? \quad (13)$$

Note that  $y$  does not appear free in  $\lambda y(x)$ , so this renaming of bound variables is legitimate. But  $\lambda x(\lambda y(x))$  is a function which, on any input  $x$ , returns the constant function that always outputs  $x$ . On the other hand,  $\lambda y(\lambda y(y))$  is a function that, on any input  $y$ , always outputs the identity function. These are different functions, and so should not be  $\alpha$ -equivalent.

To resolve this problem, we must examine the nature of substitution more carefully. This will be the task of this section.

Let us start by considering the  $\lambda$ -term  $M = \lambda x(xy)$ . This has  $y$  as a free variable. Let us consider what happens when we substitute  $x$  for  $y$ ; that is, let us calculate  $M[x/y]$ . Proceeding naively and replacing every free occurrence of  $y$  with  $x$ , it looks like we should get  $M[x/y] = \lambda x(xx)$ . This is a closed  $\lambda$ -term representing a function that takes a function  $x$  to  $xx$  - i.e., the application of  $x$  to itself.

However,  $M = \lambda x(xy)$  is the same term as  $M' = \lambda z(zy)$ , as these terms are  $\alpha$ -equivalent. But note that when we substitute  $x$  for the free variable  $y$  in  $M'$ , we find that  $M'[x/y]$  is  $\lambda z(zx)$ . This is not even a closed  $\lambda$ -term. We are left with the very awkward situation in which terms  $M$  and  $M'$  that we want to regard as the same (because they are  $\alpha$ -equivalent) return different values when the same thing (namely,  $x$ ) is substituted for the same free variable (namely,  $y$ ) in each. This is unpleasant. How should we resolve this problem?

Note that there is something peculiar about the expression  $\lambda x(xy)[x/y]$  insofar as  $x$  plays two different roles in it - first, it is something that we are substituting for  $y$ , and second, it is a variable that is already bound by  $\lambda x$  in  $\lambda x(xy)$ . Because of this, the variable  $x$  that we substitute, while free on its own, is 'captured' and becomes a bound variable in the final expression  $\lambda x(xx)$ . More generally, when evaluating  $\lambda x(xy)[E/y]$  in a case in which  $E$  contains  $x$  as a free variable, when  $E$  is substituted for  $y$  in  $\lambda x(xy)$  the free variable  $x$  in  $E$  is 'captured' by the  $\lambda x$  and becomes a bound variable. This seems undesirable, and is the source of the unpleasantness described in the previous paragraph. The natural thing to do is to prohibit this, and insist that in evaluating  $\lambda x(xy)[E/y]$  we substitute  $E$  for all free occurrences of  $y$  in a representation of  $\lambda x(xy)$  in which this sort of capturing does not happen. If for example we rewrite  $\lambda x(xy)$  as the  $\alpha$ -equivalent term  $\lambda z(zy)$  and *then* do the substitution in question, we get  $\lambda z(zy)[x/y]$ , which becomes  $\lambda z(zx)$ , and  $x$  remains free.

More generally, we will say that when evaluating  $N[E/y]$  where  $y$  is a free variable in  $N$ , we must substitute  $E$  for  $y$  in a term  $N'$   $\alpha$ -equivalent to  $N$  such that no 'capturing' of free variables occurs. In particular, we must choose some  $N'$  which is  $\alpha$ -equivalent to  $N$  such that no free variable of  $E$  becomes bound when  $E$  is substituted for every free occurrence of  $y$  in  $N'$ . The term  $N[E/y]$  is then the result of substituting  $E$  for every free occurrence of  $y$  in such an  $N'$ . Even though different choices of  $N'$  are possible, it is easily seen that the final results are all  $\alpha$ -equivalent. So for example, we could have chosen  $\lambda w(wy)$  instead of  $\lambda z(zy)$  for  $N'$  in our above example, in which case  $\lambda x(xy)[x/y]$  would become  $\lambda w(wx)$  instead of  $\lambda z(zx)$ . That is fine, as  $\lambda w(wx)$  and  $\lambda z(zx)$  are  $\alpha$ -equivalent, and are thus merely different representations of the same term.

**Definition 1.7: Capturing.**

When performing the substitution  $N[E/y]$ , a free variable  $x$  occurring in  $E$  is *captured* if some free occurrence of  $y$  in  $N$  occurs inside the scope of a  $\lambda x$  operator – that is, if  $N$  has a subexpression of the form  $\lambda x(\dots y \dots)$ , where this occurrence of  $y$  is free in  $N$ .

**Definition 1.8: Substitution.**

The  $\lambda$ -term  $N[E/y]$  is the term  $N'$  with every free occurrence of  $y$  replaced by  $E$ , where  $N'$  is a  $\lambda$ -term  $\alpha$ -equivalent to  $N$  such that no free variable of  $E$  is captured during this substitution.

That such an appropriate term  $N'$  always exists is something you will demonstrate in the exercises.

Of course, in many cases where there is no risk of capturing we may calculate  $N[E/y]$  naively by substituting  $E$  for  $y$  without carefully finding an appropriate  $N'$ . It is nevertheless important to know how to proceed when this is not the case.

Let us return now to (13). Note that in the term  $\lambda y(x)[y/x]$ , if we naively substitute  $y$  for every free occurrence of  $x$ , we get  $\lambda y(y)$ , and the variable  $y$  has been captured. Replacing  $\lambda y(x)$  with the  $\alpha$ -equivalent term  $\lambda z(x)$ , the substitution  $\lambda z(x)[y/x]$  gives the result  $\lambda z(y)$ , and there is no capturing. So  $\lambda y(x)[y/x]$  is  $\lambda z(y)$  (up to  $\alpha$ -equivalence.) Thus, equation (13) is incorrect, and must be replaced by

$$\lambda x(\lambda y(x)) =_{\alpha} \lambda y(\lambda y(x)[y/x]) = \lambda y(\lambda z(y)) \quad (14)$$

which is unproblematic. A correct understanding of substitution therefore saves our definition of  $\alpha$ -equivalence from the unpleasant example in question.

Another unusual case worth discussing is that of terms in which the same variable is bound multiple times. Consider for example the term

$$\lambda x(x \lambda x(x y)) \quad (15)$$

$\uparrow$   
 $*$

This term contains  $\lambda x$  twice. Consider the  $x$  labelled by  $*$  in the above expression. Which  $\lambda x$  should we regard this as bound by? By applying Definition 1.6, we see that (15) is  $\alpha$ -equivalent to

$$\lambda x(x \lambda z(z y)),$$

and the  $x$  in question has become a  $z$ . This shows us that in (15), we must regard the  $x$  in question as bound by the rightmost  $\lambda x$ , and not the leftmost  $\lambda x$ . (Note that we *cannot* apply Definition 1.6 to argue that (15) is  $\alpha$ -equivalent to  $\lambda z(z \lambda x(z y))$ , and so there is no analogous argument for regarding the  $x$  in question as bound by the leftmost  $\lambda x$ .) In general then, any bound variable  $x$  which occurs in the scope of multiple  $\lambda x$  abstractions will be bound by the  $\lambda x$  that occurs closest

to the left of it. Note then that on input  $c$ , the function (15) outputs the value  $(x\lambda x(xy))[c/x]$ , which is  $c\lambda x(xy)$ , and not  $c\lambda x(cy)$  or  $x\lambda x(cy)$ .

Finally, while a more intuitive grasp of substitution suffices for the ordinary manipulations of  $\lambda$ -terms encountered in much of this book, on occasions it is useful to work with a more rigorous definition of substitution, which we now present. This can safely be skipped on a first reading. In this rigorous definition,  $N[E/v]$  is defined inductively as follows:

- (i)  $v[E/v]$  is  $E$
- (ii) if  $x$  is a variable distinct from  $v$ , then  $x[E/v]$  is  $x$
- (iii)  $N(N')[E/v]$  is  $N[E/v](N'[E/v])$
- (iv)  $\lambda v(N)[E/v]$  is  $\lambda v(N)$
- (v) if  $x$  is a variable distinct from  $v$ , and either  $x$  does not occur free in  $E$  or  $v$  does not appear free in  $N$ , then  $\lambda x(N)[E/v]$  is  $\lambda x(N[E/v])$
- (vi) if  $x$  is a variable distinct from  $v$ , and  $x$  occurs free in  $E$  and  $v$  occurs free in  $N$ , then  $\lambda x(N)[E/v]$  is  $\lambda z(N[z/x][E/v])$ , where  $z$  is a variable that does not appear free or bound in  $N$  or  $E$ .

Of course, this definition only generates one of the set of  $\alpha$ -equivalent terms representing  $N[E/v]$ .

Note that by relabeling the bound variable  $x$  of  $\lambda x(N)$ , case (vi) can always be avoided. However, if it seems circular to make use of the notion of replacement of bound variables in the definition of substitution, clause (vi) may be included, and no reference to the notion of relabeling of bound variables is then necessary.

As an example, we evaluate  $\lambda r(r(\lambda s(t)))[r/t]$  using these rules.

Starting with	$\lambda r(r(\lambda s(t)))[r/t]$
Applying rule (vi), we get	$\lambda z(r(\lambda s(t))[z/r][r/t])$
Applying rule (iii), we get	$\lambda z(r[z/r](\lambda s(t)[z/r])[r/t])$
Applying rules (i) and (v), we get	$\lambda z(z(\lambda s(t[z/r]))[r/t])$
Applying rule (ii), we get	$\lambda z(z(\lambda s(t))[r/t])$
Applying rule (iii), we get	$\lambda z(z[r/t](\lambda s(t)[r/t]))$
Applying rule (ii) and (v), we get	$\lambda z(z(\lambda s(t[r/t])))$
Applying rule (i), we get	$\lambda z(z(\lambda s(r)))$

Thus  $\lambda r(r(\lambda s(t)))[r/t]$  evaluates to  $\lambda z(z(\lambda s(r)))$ , as expected.

Using this definition of substitution, facts such as the following can be rigorously proven:

**Theorem 1.9: Substitution Theorem.**

If  $x$  and  $y$  are distinct variables and  $x$  does not occur free in  $P$ , then

$$N[M/x][P/y] \text{ is identical with } N[P/y][M[P/y]/x]$$

The proof is left to the exercises.

**Exercises for Section 1.5**

1. Simplify the following expressions as far as possible, making sure to respect the requirement that free variables not be captured in substitutions. Here  $x, y$  and  $z$  are to be understood as distinct variables. (You do not have to evaluate these expressions using the more rigorous definition of substitution presented at the end of the section; instead argue more informally.)

(i)  $x(x)[y/x]$

(ii)  $y(x)[y/x]$

(iii)  $\lambda x(y)[y/x]$

(iv)  $\lambda x(yx)[y/x]$

(v)  $\lambda x(yx)[x/y]$

(vi)  $x\lambda x(yx)[x/y]$

(vii)  $x(\lambda y(xy))[\lambda x(xy)/x]$

(viii)  $x(\lambda y(xy))[\lambda y(xy)/x]$

2. Given any  $\lambda$ -terms  $N, E$  and a variable  $y$ , show that there is a term  $N'$   $\alpha$ -equivalent to  $N$  such that no free variable of  $E$  is captured in the substitution  $N'[E/y]$ .

3. Shown that  $=_\alpha$  is an equivalence relation. That is, show that for any  $\lambda$ -terms  $x, y$  and  $z$ , (i)  $x =_\alpha x$ , (ii) if  $x =_\alpha y$  then  $y =_\alpha x$ , and (iii) if  $x =_\alpha y$  and  $y =_\alpha z$ , then  $x =_\alpha z$ .

4. Perform the substitutions of problem 1. (iv), (v), (vi) and (vii) using the rigorous definition of substitution presented at the end of the section.

5. Prove Theorem 1.9. (Hint: use induction on the construction on  $N$ . When considering the case in which  $N$  has the form  $\lambda z(\bar{N})$ , you may assume by relabelling of bound variables that  $z$  is distinct from  $x$  and  $y$ , and that  $z$  does not appear free in  $M$  or  $P$ .)

6. Show that if  $x$  and  $y$  are distinct variables and  $x$  is allowed to occur free in  $P$ , then we do not necessarily have that  $N[M/x][P/y]$  is identical with  $N[P/y][M[P/y]/x]$ .

**1.6  $\beta$ -reduction and Equivalence**

In the  $\lambda$ -calculus,  $\lambda x(M)$  is the expression  $M$  viewed as a function of  $x$ . This means that it takes the input  $N$  to the output  $M[N/x]$ . But the application of the function  $\lambda x(M)$  to  $N$  can also be written as the expression  $\lambda x(M)(N)$ . The expression  $\lambda x(M)(N)$  should therefore simplify (in some sense) to  $M[N/x]$ . We call the act of simplifying

$$\lambda x(M)(N)$$

to

$$M[N/x]$$



$\beta$ -reduction, and say that  $\lambda x(M)(N)$   $\beta$ -reduces to  $M[N/x]$ . The expression  $\lambda x(M)(N)$  is sometimes called a  $\beta$ -redex, or just a redex (i.e., something that can be reduced.) The act of  $\beta$ -reducing a redex is really nothing more than the act of evaluating a function on a given input.

For example, consider the term

$$\lambda x(x(x))(\lambda z(z)).$$

This has the form  $\lambda x(M)(N)$  with  $M$  the expression  $x(x)$  and  $N$  the expression  $\lambda z(z)$ . It therefore  $\beta$ -reduces to  $x(x)[\lambda z(z)/x]$ , which is just  $\lambda z(z)(\lambda z(z))$ . This resulting term also has the form  $\lambda z(M)(N)$  with  $M$  the expression  $z$  and  $N$  the expression  $\lambda z(z)$ . It therefore  $\beta$ -reduces to  $z[\lambda z(z)/z]$ , which is just  $\lambda z(z)$ . When  $A$   $\beta$ -reduces to  $B$  and  $B$   $\beta$ -reduces to  $C$ , we also say that  $A$   $\beta$ -reduces to  $C$ . Thus  $\lambda x(x(x))(\lambda z(z))$   $\beta$ -reduces to  $\lambda z(z)$ .

We will allow  $\beta$ -reduction to be performed on a redex inside a term. Thus, the term

$$\tau_1 \lambda x(M)(N) \tau_2$$

$\beta$ -reduces to

$$\tau_1 M[N/x] \tau_2.$$

where  $\tau_1$  and  $\tau_2$  are sequences of symbols. For example, consider the term

$$\lambda y(\lambda x(x(x))(\lambda z(z)))(y)$$

This has the form  $\tau_1 \lambda x(M)(N) \tau_2$  with  $\tau_1$  the sequence of symbols  $\lambda y(,$  and  $M$  the expression  $x(x)$ , and  $N$  the expression  $\lambda z(z)$ , and  $\tau_2$  the sequence of symbols  $(y))$  – i.e.,

$$\underbrace{\lambda y(}_{\tau_1} \underbrace{\lambda x(x(x))}_{M} \underbrace{(\lambda z(z))}_{N} \underbrace{(y))}_{\tau_2}$$

It therefore  $\beta$ -reduces to

$$\lambda y(x(x)[\lambda z(z)/x])(y)$$

which is

$$\lambda y(\lambda z(z)(\lambda z(z)))(y).$$

This in turn has the form  $\tau_1 \lambda z(M)(N) \tau_2$  with  $\tau_1$  the expression  $\lambda y($  and  $M$  the expression  $z$  and  $N$  the expression  $\lambda z(z)$  and  $\tau_2$  the expression  $(y))$  – i.e.,

$$\underbrace{\lambda y(}_{\tau_1} \underbrace{\lambda z(z)}_{M} \underbrace{(\lambda z(z))}_{N} \underbrace{(y))}_{\tau_2}$$

It therefore  $\beta$ -reduces to

$$\lambda y(z[\lambda z(z)/z])(y)$$

which is

$$\lambda y(\lambda z(z)(y)).$$

You should be able to see that by similar reasoning, this  $\beta$ -reduces to  $\lambda y(y)$ , which cannot be  $\beta$ -reduced any further. Thus the original expression  $\lambda y(\lambda x(x(x))(\lambda z(z))(y))$   $\beta$ -reduces to  $\lambda y(y)$ .

One must be careful here and pay very close attention to the presence or absence of parentheses. In order to apply  $\beta$ -reduction within a term, the overall term just have the *exact* form

$$\tau_1 \lambda x(M)(N) \tau_2.$$

In particular, there must a set of parentheses immediately surrounding both  $M$  and  $N$  as depicted, and *nothing* between  $(M)$  and  $(N)$ . With this is mind, consider the exaxmple:

$$y(\lambda x(x))(z).$$

We *cannot* apply the function  $\lambda x(x)$  to  $z$  and  $\beta$ -reduce the whole expression to  $yz$ , letting  $\tau_1$  be  $y$  and  $M$  be the expression  $x$  and  $N$  be the expression  $z$  and  $\tau_2$  the empty expression – i.e.,

$$\underbrace{\lambda y(\lambda x(\underbrace{x}_M))(\underbrace{z}_N))}_{\tau_1}$$

because  $(M)$  does not sit *immediately* next to  $(N)$ , due to the intervening  $)$ . The  $\lambda$ -term  $y(\lambda x(x))(z)$  in fact contains no redexes, and thus cannot be  $\beta$ -reduced in any way.

Determining whether it is correct to perform a  $\beta$ -reduction can sometimes be confusing when have used our conventions to abbreviate expressions and remove parentheses. To check whether one can apply  $\beta$ -reduction within a term, it is often best to restore any missing parentheses first, just to be sure. Paying attention to possibly suppressed parentheses is especially important also because of the failure of associativity of the  $\lambda$ -calculus: the  $\lambda$ -terms  $x(yz)$  and  $(xy)z$  are not in general the same in the  $\lambda$ -calculus, and so the placement of parentheses can completely change the meaning of a  $\lambda$ -term. You will explore this point in the exercises.

When a  $\lambda$ -term  $A$  reduces to  $B$  by a single  $\beta$ -reduction, we write  $A \rightarrow_\beta B$ . When  $A$  reduces to  $B$  by a *sequence* of  $\beta$ -reductions (this includes the possibility that the sequence consists of just a single  $\beta$ -reduction), we write  $A \twoheadrightarrow_\beta B$ . So we have just seen that

$$\lambda y(\lambda x(x(x))(\lambda z(z))(y)) \twoheadrightarrow_\beta \lambda y(y).$$

We also write  $A \twoheadrightarrow_\beta A$  for any term  $A$ , thinking of there as being a ‘0 step’  $\beta$ -reduction of  $A$  to  $A$ .

We summarize these definitions as follows.

**Definition 1.10:  $\beta$ -reduction.**

An expression of the form  $\lambda x(M)(N)$  (where  $x$  is a variable and  $M, N$  are  $\lambda$ -terms) is called a  $\beta$ -redex (or just a redex). We have:

- (i)  $A \rightarrow_\beta B$  iff  $A$  has the form  $\tau_1 \lambda x(M)(N) \tau_2$  and  $B$  has the form  $\tau_1 M[N/x] \tau_2$ .
- (ii)  $A \twoheadrightarrow_\beta B$  ( $A$   $\beta$ -reduces to  $B$ ) iff for some sequence  $X_0, X_1, \dots, X_n$  of  $\lambda$ -terms we have that

$A = X_0, B = X_n$ , and  $X_0 \rightarrow_{\beta} X_1 \rightarrow_{\beta} \dots \rightarrow_{\beta} X_n$ . This includes the degenerate case in which  $A$  and  $B$  are identical.

The set of rules of formation of  $\lambda$ -terms along with the rule of  $\beta$ -reduction just given can be thought of as a system of rules for manipulating terms constructed in a certain sort of way. This system is typically called the *untyped  $\lambda$ -calculus*. We will also call it  $\lambda_0$ , to easily distinguish it from expansions of the untyped  $\lambda$ -calculus that we consider later.

**Definition 1.11: Rules of the Untyped  $\lambda$ -calculus  $\lambda_0$ .**

Rules of  $\lambda$ -term formation

- (i) Any variable  $x, y, z, \dots$  is a  $\lambda$ -term.
- (ii) If  $M$  and  $N$  are  $\lambda$ -terms, then so is  $M(N)$ .
- (iii) If  $M$  is a  $\lambda$ -term and  $x$  a variable, then  $\lambda x(M)$  is a  $\lambda$ -term.

Rule of  $\beta$ -reduction

- (i)  $\lambda x(M)(N) \rightarrow_{\beta} M[N/x]$

Note that for compactness, in the above definition we state only the simplest form of the rule of  $\beta$ -reduction, with it being understood that  $\beta$ -reduction may be performed on redexes inside terms.

We now have the following important definition

**Definition 1.12:  $\beta$ -normal form.**

A term  $X$  is in  $\beta$ -normal form just in case no  $\beta$ -reduction may be performed on it – that is, just in case it contains no subexpression of the form  $\lambda x(M)(N)$ .

For example,  $\lambda y(y)$  is in  $\beta$ -normal form. By performing a sequence of  $\beta$ -reductions, we reduced  $\lambda y(\lambda x(x(x))(\lambda z(z))(y))$  to the  $\beta$ -normal form term  $\lambda y(y)$ . One might hope that starting with any term and continually  $\beta$ -reducing it, one will eventually arrive at a term in  $\beta$ -normal form. Unfortunately this is not the case. For example, consider the term

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

We have

$$\begin{aligned} \lambda x(x(x))(\lambda x(x(x))) &\rightarrow_{\beta} x(x)[\lambda x(x(x))/x] \\ &= \lambda x(x(x))(\lambda x(x(x))) \end{aligned}$$

Although the term  $\lambda x(x(x))(\lambda x(x(x)))$  is not in  $\beta$ -normal form, the only act of  $\beta$ -reduction that may be performed on it does not simplify the original term, but rather just returns the term itself.

There is thus no sequence of  $\beta$ -reductions that begins with this term and ends with a term in  $\beta$ -normal form.

A yet further complication that needs to be considered is that there are  $\lambda$ -terms on which there are *multiple*  $\beta$ -reductions that can be performed, because they contain multiple redexes. For example, consider the term

$$\lambda x(x)(y(y))(\lambda z(z)(y))$$

If we perform a  $\beta$ -reduction on the redex  $\lambda x(x)(y(y))$ , this redex becomes  $y(y)$ , and so the whole term  $\beta$ -reduces to  $y(y)(\lambda z(z)(y))$ . This in turn  $\beta$ -reduces to  $y(y)(y)$  (i.e.,  $yyy$ ), which is in  $\beta$ -normal form.

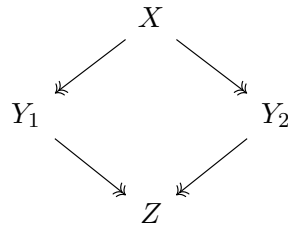
However, we can also perform a  $\beta$ -reduction on the redex  $\lambda z(z)(y)$  first. Doing so, the entire term is  $\beta$ -reduced to  $\lambda x(x)(y(y))(y)$ , which in turn  $\beta$ -reduces to  $y(y)(y)$  (i.e.  $yyy$ .)

In this case, even though there are different ‘orders’ in which we can  $\beta$ -reduce our given term, different  $\beta$ -reductions ultimately give the same  $\beta$ -normal term. Is this true in general? It turns out that it is. This is a consequence of the following important general fact:

**Theorem 1.13: The Church-Rosser Theorem for  $\beta$ -reductions in  $\lambda_0$ .**

In the untyped  $\lambda$ -calculus  $\lambda_0$ , if  $X \rightarrow_{\beta} Y_1$  and  $X \rightarrow_{\beta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ .

This theorem is often depicted as follows:



This theorem tells us that if different sequences of  $\beta$ -reductions take a term  $X$  to different terms  $Y_1$ ,  $Y_2$ , then *further*  $\beta$ -reductions can take  $Y_1$  and  $Y_2$  back to a common term  $Z$ . So any distinct sequences of  $\beta$ -reductions can be ‘reunited’, forming the diamond shape shown above.

It follows from this that when a term  $X$  can be reduced to a  $\beta$ -normal form  $Y$ , this  $Y$  is unique (up to  $\alpha$ -equivalence). That is to say, it follows that different chains of  $\beta$ -reductions cannot take us from an  $X$  to two distinct  $\beta$ -normal forms  $Y_1$  and  $Y_2$ . To see why, suppose to the contrary that two different chains of  $\beta$ -reductions take us from a term  $X$  to two  $\beta$ -normal forms  $Y_1$  and  $Y_2$ . It follows from the Church-Rosser Theorem just stated there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ . But if  $Y_1$  and  $Y_2$  are in  $\beta$ -normal form then they cannot be reduced any further, and so we must have that  $Y_1$  and  $Z$  are the same term (up to  $\alpha$ -equivalence) and  $Y_2$  and  $Z$  are the same term (up to  $\alpha$ -equivalence). By transitivity of  $\alpha$ -equivalence, it follows that  $Y_1$  and  $Y_2$  are also  $\alpha$ -equivalent.

So if different chains of  $\beta$ -reductions take us from  $X$  to distinct  $\beta$ -normal forms  $Y_1$  and  $Y_2$ , then  $Y_1$  and  $Y_2$  are  $\alpha$ -equivalent.

The proof of the Church-Rosser Theorem is perhaps surprisingly complex. We include it in the appendix to this chapter.

Intuitively,  $\beta$ -reduction should be thought of as a process taking a  $\lambda$ -term into another  $\lambda$ -term referring to the same entity. Consider then the following ‘symmetrical’ version of one-step  $\beta$ -reduction:

**Definition 1.14: The relation  $\sim_\beta$ .**

$X \sim_\beta Y$  if either  $X \rightarrow_\beta Y$  or  $Y \rightarrow_\beta X$

We have  $X \sim_\beta Y$  iff there is a one-step  $\beta$ -reduction from one of  $X$  or  $Y$  to the other. If  $X \sim_\beta Y$ , then  $X$  and  $Y$  should refer to the same entity. Because ‘referring to the same entity’ is surely a transitive relation, we should then have that if  $X_0 \sim_\beta X_1, X_1 \sim_\beta X_2, \dots, X_{n-1} \sim_\beta X_n$ , then  $X_0$  and  $X_n$  ‘refer to the same entity’. With this in mind, we define the relation of  $\beta$ -equivalence as follows:

**Definition 1.15: The relation  $U =_\beta V$ .**

$U$  and  $V$  are  $\beta$ -equivalent (in symbols,  $U =_\beta V$ ) iff one of the following holds:

- (a)  $U$  and  $V$  are the same term (up to  $\alpha$ -equivalence), or
- (b) for some sequence of terms  $X_0, X_1, \dots, X_{n-1}, X_n$  with  $U = X_0$  and  $V = X_n$ ,

$$X_0 \sim_\beta X_1, X_1 \sim_\beta X_2, \dots, X_{n-1} \sim_\beta X_n.$$

Although when  $U$  and  $V$  are  $\beta$ -equivalent we regard them as referring to the same thing, this does *not* mean that we regard them as the same expression. (This is a common feature of any language. In English the terms ‘Prince Charles’ mother’ and ‘Queen Elizabeth II’ refer to the same person, even though they are different expressions.) In this sense, what is being said here is different from what was said about  $\alpha$ -equivalence, where certain distinct lexical items were regarded as literally the same expression. We do not want to make that move here, as we want it to be possible for syntactically distinct terms to have a common referent.

The following is a useful corollary of the Church-Rosser Theorem

**Corollary 1.16**

In  $\lambda_0$ ,  $U =_\beta V$  iff there is a term  $Z$  such that  $U \rightarrow_\beta Z$  and  $V \rightarrow_\beta Z$ .

You will prove this in an exercise. This corollary tells us that two  $\lambda$ -terms are  $\beta$ -equivalent iff they can be  $\beta$ -reduced to a common term.

A final fact about  $\beta$ -reduction that is occasionally useful is the following

**Lemma 1.17:**

For any  $\lambda_0$ -terms  $U, V, t$  and variable  $x$ , if  $U \rightarrow_\beta V$  then  $U[t/x] \rightarrow_\beta V[t/x]$

The proof is also left to the exercises.

**Exercises for Section 1.6**

1. Reduce each of the following lambda terms into  $\beta$ -normal form. In all cases, if an expression is already in  $\beta$ -normal form, you simply need to point this out and do nothing further. Here,  $x, y, z, u$  and  $v$  are to be understood as distinct variables.

- (i).  $\lambda x(x)(y)$
- (ii).  $\lambda x(E)(z)$  (Here  $E$  is a term that contains no free occurrence of  $x$ .)
- (iii).  $\lambda x(x)(\lambda y(y(y)))$
- (iv).  $\lambda x(x)(\lambda y(y(y)))(\lambda z(z))$
- (v).  $\lambda x(y(\lambda y(z)))$
- (vi).  $\lambda x(\lambda y(xy))(u)$
- (vii).  $\lambda x(\lambda y(xy))(\lambda u(u(v)))$
- (viii).  $\lambda x(\lambda y(xy))(\lambda u(\lambda v(u(v))))$

2. In this exercise, we show that the application of  $\lambda$ -terms is not associative. In particular, we show that it is *not* the case that for all  $\lambda$ -terms  $a, b$  and  $c$ ,  $a(b)(c)$  and  $a(b(c))$  are  $\beta$ -equivalent. Let  $x$  and  $y$  be distinct variables, and let  $a$  be  $\lambda x(y)$ ,  $b$  be  $y$ , and  $c$  be  $y$ . Show that  $a(b)(c)$   $\beta$ -reduces to  $y(y)$  and  $a(b(c))$   $\beta$ -reduces to  $y$ . Because both  $y$  and  $y(y)$  are in  $\beta$ -normal form and are different terms (i.e., not  $\alpha$  equivalent), it follows that  $y$  and  $y(y)$  are not  $\beta$ -equivalent terms.

3. In spite of the previous exercise, the  $\lambda$ -calculus captures the fact that function composition is associative. Given functions (i.e.,  $\lambda$ -terms)  $f$  and  $g$ , define  $[f \circ g]$  to be  $\lambda x(f(g(x)))$ . Show that the terms  $[f \circ g](h(x))$  and  $f([g \circ h](x))$  are  $\beta$ -equivalent. How are the results of this and the previous exercise compatible? Doesn't one exercise show that function composition is associative, and the other show that it is not?

4. Argue that if  $s$  is a closed  $\lambda$ -term and  $s \rightarrow_{\beta} t$ , then  $t$  is a closed  $\lambda$ -term. (A closed term is one with no free variables.)
5. Argue that every  $\beta$ -normal closed term can be written in the form

$$\lambda x_1 \dots \lambda x_n (x(M_1) \dots (M_m))$$

where  $n > 0, m \geq 0$ ,  $x$  and the  $x_i$  are (not necessarily distinct) variables, and the  $M_i$  are all in  $\beta$ -normal form. (Hint: prove this by induction on the construction of  $\lambda$ -terms.)

6. Using the Church-Rosser Theorem, prove Corollary 1.16.
7. Show that there is a  $\lambda$ -term  $\tau$  which has a non-terminating  $\beta$ -reduction sequence - i.e., an infinite sequence  $\tau \rightarrow_{\beta} \tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$  - as well as a terminating  $\beta$ -reduction sequence - i.e., a finite sequence  $\tau \rightarrow_{\beta} \sigma_1 \rightarrow_{\beta} \sigma_2 \rightarrow_{\beta} \dots \rightarrow_{\beta} \sigma_i$  for some finite  $i$ , where  $\sigma_i$  is in  $\beta$ -normal form.
8. Find terms  $X$  and  $Y$  such that neither  $X$  nor  $Y$  has a  $\beta$ -normal form, but  $X(Y)$  has a  $\beta$ -normal form.
9. Let  $Y$  be the term  $\lambda x (\lambda y (x(yy)) (\lambda y (x(yy))))$ . For every term  $E$ , show that

$$E(Y(E)) =_{\beta} Y(E).$$

Thus,  $Y(E)$  is a 'fixed point' of  $E$ .

10. Show that the term  $Y$  of the previous problem does not have a  $\beta$ -normal form. Show that in fact for any term  $Y'$  such that for every term  $E$ ,  $E(Y'(E)) =_{\beta} Y'(E)$ ,  $Y'$  does not have a  $\beta$ -normal form. (Hint: let  $x$  be a variable that does not appear free in  $Y'$ , and consider the fact that  $x(Y'(x)) =_{\beta} Y'(x)$ .)
11. Prove Lemma 1.17. (Hint: by  $\alpha$ -equivalence, you may assume that neither  $x$  nor any free variable of  $t$  is bound in  $V$ .)
12. Is it true that for any  $\lambda$ -terms  $U, V, t$  and variable  $x$ , if  $U \rightarrow_{\beta} V$  then  $U[t/x] \rightarrow_{\beta} V[t/x]$ ?

## 1.7 $\eta$ -reduction and Extensionality

Suppose  $E$  is a term in the  $\lambda$ -calculus that does not contain  $x$  as a free variable. Consider the term  $E^* = \lambda x (E(x))$ . On any input  $c$ , we have

$$\begin{aligned} E^*(c) &= \lambda x (E(x))(c) \\ &=_{\beta} E(x)[c/x] \\ &= E(c) \end{aligned}$$

As functions,  $E^*$  and  $E$  produce the same outputs on the same inputs. It is thus tempting to think of them as the same function.

Note that if we do not require that  $E$  does not contain  $x$  as a free variable, then  $E^*(c)$  is not necessarily  $\beta$ -equivalent to  $E(c)$ . To see this, let  $x$  and  $y$  be distinct variables, and let  $E$  just be the variable  $x$ . Then  $E(y)$  is  $x(y)$ , but  $E^*(y)$  is  $\lambda x(x(x))(y)$ , which  $\beta$ -reduces to  $y(y)$ , and so  $E(y)$  and  $E^*(y)$  are not  $\beta$ -equivalent.

Suppose then that  $E$  does not contain  $x$  as a free variable, so that  $E^*$  and  $E$  produce the same outputs on the same inputs. In spite of this, nothing we have said so far forces us to identify  $E^*$  and  $E$  in any sense. Consider for example the special case in which  $E$  is just the variable  $y$ . In this case,  $E^*$  is  $\lambda x(yx)$ . Both  $y$  and  $\lambda x(yx)$  are terms in  $\beta$ -normal form, and are not  $\alpha$ -equivalent. Thus, they are not  $\beta$ -equivalent either. So  $E^*$  and  $E$  are not equivalent in any formal sense we have identified yet.

In the  $\lambda$ -calculus as we have developed it thus far, we do not necessarily treat a pair of functions as identical just because they always give the same outputs on the same inputs. Is this a defect or a virtue? Consider the following principle

**Definition 1.18: Principle of Extensionality.**

If two terms  $\lambda$ -terms  $U$  and  $V$  have the property that for every  $\lambda$ -term  $c$ ,  $U(c)$  and  $V(c)$  are the same, then  $U$  and  $V$  are the same.

This principle is not fully precise, insofar as the term ‘same’ has not yet been defined. However, if by ‘same’ we mean  $\beta$ -equivalent, then we have seen that this principle fails.

Perhaps this is not worrisome. One might, after all, have reservations about the Principle of Extensionality. Must we really identify two functions or algorithms that each return the same output on the same inputs? Two algorithms that always return the same output on the same inputs might still be different in many other respects - one algorithm might require significantly more time or memory to produce an output than the other, for example. Whether this sort of distinction matters will depend on our purposes. If we are interested in implementations of functions on real physical computers, such distinctions may well be important. If by contrast we are only interested in functions from a very abstract mathematical point of view, such distinctions may well not matter. There is no right or wrong answer as to whether we must endorse the Principle of Extensionality. It simply depends on whether doing so is useful to our ends, whatever those are.

Suppose however that we are attracted to the Principle of Extensionality, and are interested in articulating a notion of sameness according to which it holds. How might we proceed? To start with, for any term  $E$  we will want to identify  $\lambda x(E(x))$  (where  $x$  is not free in  $E$ ) with  $E$ . To this end, we introduce a new reduction rule - the rule of  $\eta$ -reduction - according to which the term  $\lambda x(E(x))$  (with  $x$  not free in  $E$ ) can be reduced to  $E$ , as follows

If  $x$  not free in  $E$ , then  $\lambda x(E(x)) \rightarrow_{\eta} E$ .

In this definition, we are of course free to use any variable in place of  $x$ . The expression  $\lambda x(E(x))$  is sometimes called an  $\eta$ -redex, or just a redex. As with  $\beta$ -reduction, we allow ourselves to apply



$\eta$ -reductions to  $\eta$ -redexes within terms. In this way, the term

$$\tau_1 \lambda x(E(x)) \tau_2$$

also  $\eta$ -reduces to

$$\tau_1 E \tau_2$$

(where  $\tau_1$  and  $\tau_2$  are sequences of symbols, and  $x$  is not free in  $E$ .)

We thus regard both  $\beta$ -reduction and  $\eta$ -reduction as valid ways of simplifying terms. In this spirit, we define  $X \rightarrow_{\beta\eta} Y$  iff  $X \rightarrow_{\beta} Y$  or  $X \rightarrow_{\eta} Y$ , and  $X \twoheadrightarrow_{\beta\eta} Y$  iff some sequence of  $\beta$  or  $\eta$  reductions takes  $X$  to  $Y$  (where again, this includes the degenerate ‘0 step’ reduction in which  $X$  and  $Y$  are  $\alpha$ -equivalent.) More formally, we have the following definitions:

**Definition 1.19:  $\eta$ -reduction,  $\beta\eta$ -reduction, and  $\beta\eta$ -equivalence**

An expression of the form  $\lambda x(E(x))$  (where  $x$  is a variable and  $E$  a  $\lambda$ -term such that  $x$  does not appear free in  $E$ ) is called an  $\eta$ -redex (or just a *redex*). We have:

(a)  $A \rightarrow_{\eta} B$  iff  $A$  has the form  $\tau_1 \lambda x(E(x)) \tau_2$  and  $B$  has the form  $\tau_1 E \tau_2$  for expressions  $\tau_1, E, \tau_2$ , and  $x$  is not free in  $E$ .

(b)  $A \rightarrow_{\beta\eta} B$  iff  $A \rightarrow_{\beta} B$  or  $A \rightarrow_{\eta} B$ .

(c)  $A \twoheadrightarrow_{\beta\eta} B$  ( $A$   $\beta\eta$ -reduces to  $B$ ) iff for some sequence  $X_0, X_1, \dots, X_n$  of  $\lambda$ -terms we have that  $A = X_0, B = X_n$ , and  $X_0 \rightarrow_{\beta\eta} X_1 \dots \rightarrow_{\beta\eta} X_n$ . This includes the degenerate case in which  $A$  and  $B$  are identical.

(d)  $X$  and  $Y$  are  $\beta\eta$ -equivalent (in symbols,  $X =_{\beta\eta} Y$ ) iff one of the following holds:

- (i)  $X$  and  $Y$  are the same term (up to  $\alpha$ -equivalence), or
- (ii) for some sequence of terms  $X_0, X_1, \dots, X_{n-1}, X_n$  with  $U = X_0$  and  $V = X_n$ ,

$$X_0 \sim_{\beta\eta} X_1, X_1 \sim_{\beta\eta} X_2, \dots, X_{n-1} \sim_{\beta\eta} X_n.$$

In fact,  $\beta\eta$ -reduction and  $\beta\eta$ -equivalence have many of the same important properties that  $\beta$ -reduction and  $\beta$ -equivalence have. For example

**Theorem 1.20: The Church-Rosser Theorem for  $\beta\eta$ -reductions in  $\lambda_0$ .**

In the untyped  $\lambda$ -calculus  $\lambda_0$ , if  $X \twoheadrightarrow_{\beta\eta} Y_1$  and  $X \twoheadrightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \twoheadrightarrow_{\beta\eta} Z$  and  $Y_2 \twoheadrightarrow_{\beta\eta} Z$ .

Again, the proof of this result is left to the appendix. As before, the following corollary is easy:

**Corollary 1.21**

In  $\lambda_0$ ,  $U =_{\beta\eta} V$  iff there is a term  $Z$  such that  $U \rightarrow_{\beta\eta} Z$  and  $V \rightarrow_{\beta\eta} Z$ .

We also have the following definition

**Definition 1.22:  $\beta\eta$ -normal form.**

A term  $X$  is in  $\beta\eta$ -normal form just in case no  $\beta$ -reduction or  $\eta$ -reduction may be performed on it – i.e., just in case it contains no subexpression of the form  $\lambda x(M)(N)$  or of the form  $\lambda x(E(x))$ , where  $x$  is not free in  $E$ .

Corollary 1.21 then shows that  $\beta\eta$ -equivalent terms can be  $\beta\eta$ -reduced to a common term. It also follows from Theorem 1.20 as before that when a term can be reduced to a term in  $\beta\eta$ -normal form, this term in  $\beta\eta$ -normal form is unique (up to  $\alpha$ -equivalence.)

In analogy with Lemma 1.17, a further useful fact about  $\eta$ -reduction is the following

**Lemma 1.23:**

For any  $\lambda_0$ -terms  $U, V, t$  and variable  $x$ , if  $U \rightarrow_{\eta} V$  then  $U[t/x] \rightarrow_{\eta} V[t/x]$

The proof is also left to the exercises.

An interesting technical property of  $\beta\eta$ -reductions is that whenever an  $\eta$ -reduction is followed by a  $\beta$ -reduction, there is a sense in which the order of these reductions can be reversed, so that the  $\beta$ -reduction comes before the  $\eta$ -reduction. More precisely, we have the following lemma

**Lemma 1.24:  $\eta$ -reduction postponement in  $\lambda_0$ .**

For all  $\lambda$ -terms  $X, Y, Z$ , if  $X \rightarrow_{\eta} Y \rightarrow_{\beta} Z$  in  $\lambda_0$ , then there is a  $\lambda$ -term  $Y'$  such that  $X \rightarrow_{\beta} Y' \rightarrow_{\eta} Z$  (where  $\rightarrow_{\eta}$  is a sequence of  $\eta$ -reductions.)

You will prove this in the exercises. As a consequence, it can be shown that if a term can be reduced to  $\beta\eta$ -normal form, then there is a way of reducing it to  $\beta\eta$ -normal form that involves performing a sequence of  $\beta$ -reductions followed by a sequence of  $\eta$ -reductions.

All things considered, the notion of  $\beta\eta$ -equivalence is a natural notion with many pleasing properties. In fact, the notion of  $\beta\eta$ -equivalence provides us with a notion of ‘sameness’ which it turns out is sufficient to demonstrate the Principle of Extensionality. In particular, we have the following

**Theorem 1.25: Principle of Extensionality for  $\beta\eta$ -equivalence.**

In the untyped  $\lambda$ -calculus  $\lambda_0$ , if two terms  $\lambda$ -terms  $U$  and  $V$  have the property that for every  $\lambda$ -term  $c$ ,  $U(c) =_{\beta\eta} V(c)$ , then  $U =_{\beta\eta} V$ .

**Proof**

Suppose we have  $\lambda$ -terms  $U$  and  $V$  such that for every  $\lambda$ -term  $c$ ,  $U(c) =_{\beta\eta} V(c)$ . Let  $z$  be a variable that does not appear free in  $U$  or  $V$ . Then  $U(z) =_{\beta\eta} V(z)$ . In an exercise, you will show that if  $E_1$  and  $E_2$  are  $\lambda$ -terms such that  $E_1 =_{\beta\eta} E_2$  and  $z$  is a variable, then  $\lambda z(E_1) =_{\beta\eta} \lambda z(E_2)$ . We therefore have that  $\lambda z(U(z)) =_{\beta\eta} \lambda z(V(z))$ . Because  $\lambda z(U(z)) =_{\eta} U$  and  $\lambda z(V(z)) =_{\eta} V$ , combining these and using the symmetry and transitivity of  $=_{\beta\eta}$  gives  $U =_{\beta\eta} V$ .  $\square$

The relation  $=_{\beta\eta}$  is therefore a notion of sameness of  $\lambda$ -terms for which the Principle of Extensionality holds.

There is room to wonder whether the Principle of Extensionality as stated in Definition 1.18 adequately captures the intuitive idea that functions which give the same outputs on the same inputs are the same function. For instance, it might be suggested that we should be willing to identify  $\lambda$ -terms  $U$  and  $V$  for which  $U(c)$  and  $V(c)$  are the same for every *closed*  $\lambda$ -term  $c$ . Or perhaps we should identify  $\lambda$ -terms  $U$  and  $V$  for which  $U(c)$  and  $V(c)$  are the same merely for every closed  $\lambda$ -term  $c$  in  $\beta$ -normal form. This would lead to various modifications of the Principle of Extensionality. Whether new reduction rules can be found that render the corresponding modified versions of Theorem 1.25 true is not clear. More generally, whether we should be content with Definition 1.18 and Theorem 1.25 as a statement of the ‘extensionality’ of our system is a question we shall have to leave open.

At this point, we have two independent notions of reduction available to us:  $\beta$ -reduction and  $\eta$ -reduction. We also have the notion of reduction that consists of their combination, namely,  $\beta\eta$ -reduction. All three of these notions of reduction satisfy the Church-Rosser Theorem, and so they have much in common. In most texts,  $\beta$ - and  $\beta\eta$ -reduction are the most closely studied. The notion of  $\beta$ -reduction is especially natural given the interpretation of  $\lambda$ -terms as functions which can operate on arguments, and the notion of  $\beta\eta$ -reduction is natural because of the way it ensures a type of extensionality. Both notions will be important in what follows. We will expand and explore these notions of reduction further in following chapters.

**Exercises for Section 1.7**

1. Reduce each of the following terms to  $\beta\eta$ -normal form, or argue that they cannot be reduced to  $\beta\eta$ -normal form. Here,  $x$  and  $y$  are distinct variables.

- (i).  $\lambda x(x(x))(\lambda x(y(x)))$
- (ii).  $\lambda x(y(x))(\lambda x(x(x)))$
- (iii).  $\lambda x(y(x))(\lambda x(y(x)))$
- (iv).  $\lambda x(x(x))(\lambda x(x(x)))$

2. Prove that if  $E_1$  and  $E_2$  are  $\lambda$ -terms such that  $E_1 =_{\beta\eta} E_2$  and  $v$  is a variable, then  $\lambda v(E_1) =_{\beta\eta} \lambda v(E_2)$ . (Hint: use induction on  $n$ , where  $X_0, X_1, \dots, X_{n-1}, X_n$  is a sequence of  $\lambda$ -terms with  $U = E_1$  and  $V = E_2$  such that  $X_0 \sim_{\beta\eta} X_1, X_1 \sim_{\beta\eta} X_2, \dots, X_{n-1} \sim_{\beta\eta} X_n$ .)

3. Prove that no  $\lambda$ -term has an infinite  $\eta$ -reduction sequence. That is to say, prove that there is no sequence of terms  $X_0, X_1, X_2, \dots$  such that  $X_0 \rightarrow_{\eta} X_1 \rightarrow_{\eta} X_2 \rightarrow_{\eta} \dots$

4. Prove Lemma 1.24. (Hint: suppose the term  $X$  contains an  $\eta$ -redex  $\lambda x(E(x))$  which is  $\eta$ -reduced to  $E$  in the expression  $Y$ . We are supposing that  $Y$  then contains a  $\beta$ -redex  $\lambda y(F)(G)$  which when  $\beta$ -reduced gives  $Y \rightarrow_{\beta} Z$ . From the Subterm Lemma (Lemma 1.5), in the term  $Y$  either the subterms  $E$  and  $\lambda y(F)(G)$  are identical, one is properly contained in the other, or they are disjoint.

We thus have the following cases to consider:

Case (i):  $E$  and  $\lambda y(F)(G)$  are identical subterms of  $Y$ . This means that  $X$  has the form

$$X = \dots \lambda x(\lambda y(F)(G)(x)) \dots$$

Case (ii): In the term  $Y$ , the subterm  $E$  is properly contained in the subterm  $\lambda y(F)(G)$ . This means that in the term  $Y$ , (a)  $E$  is the term  $\lambda y(F)$ , (b)  $E$  is contained in  $F$ , or (c)  $E$  is contained in  $G$ . In subcase (a), this means that  $X$  has the form

$$X = \dots \lambda x(\lambda y(F)(x))(G) \dots$$

In subcase (b) this means that  $Y$  has the form

$$Y = \dots \lambda y(\dots E \dots)(G) \dots$$

and so  $X$  has the form

$$X = \dots \lambda y(\dots \lambda x(E(x)) \dots)(G) \dots$$

and in subcase (c), this means that  $X$  has the form

$$X = \dots \lambda y(F)(\dots \lambda x(E(x)) \dots) \dots$$

Case (iii): In the term  $Y$ , the subterm  $\lambda y(F)(G)$  is properly contained in the subterm  $E$ . This means that  $X$  has the form

$$X = \dots \lambda x(\dots \lambda y(F)(G) \dots)(x) \dots$$

Case (iv):  $E$  and  $\lambda y(F)(G)$  are disjoint subterms of  $Y$ . Then it is easily argued that  $\lambda x(E(x))$  must also be disjoint from  $\lambda y(F)(G)$  in  $X$ , and so  $X$  has the form

$$X = \dots \lambda x(E(x)) \dots \lambda y(F)(G) \dots$$

Consider all these possible forms for  $X$  and show that the result holds in each case.)

5. Using Lemma 1.24, prove the following generalization of Lemma 1.24

For all  $\lambda$ -terms  $X$  and  $Z$ , if  $X \rightarrow_{\beta\eta} Z$ , then there is a  $\lambda$ -term  $Y$  such that  $X \rightarrow_{\beta} Y \rightarrow_{\eta} Z$ .

6. Prove that if a term has an infinite  $\beta\eta$ -reduction sequence, then it also has an infinite  $\beta$ -reduction sequence. (Hint: use the result of the previous exercise.)

7. (i) Show that if  $X \rightarrow_{\eta} Y$  and  $X$  is in  $\beta$ -normal form, then  $Y$  is in  $\beta$ -normal form. (ii) Give an example of  $X \rightarrow_{\eta} Y$  in which  $Y$  is in  $\beta$ -normal form, but  $X$  is not in  $\beta$ -normal form.

8. In contrast with (ii) of the previous exercise, show that if  $X \rightarrow_{\eta} Y$  and  $Y$  is in  $\beta$ -normal form, then  $X$  has a  $\beta$ -normal form.

9. Prove that if a term has a  $\beta$ -normal form, then it has a  $\beta\eta$ -normal form.

10. Prove that if a term has a  $\beta\eta$ -normal form, then it has a  $\beta$ -normal form.

11. Prove Lemma 1.23.

## 1.8 The untyped $\lambda$ -calculus as a model of computation - I

By now you should have some confidence with the manipulation of  $\lambda$ -terms. Still, at this point all these manipulations might seem like pure formalism.  $\lambda$ -terms themselves are functions, but they are very strange sorts of functions that take *any* other function as a possible input, and so are quite different from the sorts of functions of everyday mathematical practice. What is then the mathematical point of this strange system?

For our purposes, one of the main points of the  $\lambda$ -calculus is that it is a very powerful *programming language*. In this section, we will show how certain traditional calculations or computations with which you are familiar can be interpreted as  $\lambda$ -calculus manipulations. In the next section, we will talk more abstractly about the  $\lambda$ -calculus as a programming language or framework for computation.

As our first example of the utility of the  $\lambda$ -calculus, let us focus on the manipulation of truth values. We begin by noting that both the terms  $\lambda x(\lambda y(x))$  and  $\lambda x(\lambda y(y))$  are in  $\beta$ -normal form and so cannot be simplified any further by  $\beta$ -reduction. Let us think of the former as denoting the truth value **T** and the later as denoting the truth value **F**:

$$\mathbf{T} \equiv \lambda x(\lambda y(x)), \quad \mathbf{F} \equiv \lambda x(\lambda y(y)). \quad (16)$$

Consider then the function  $\lambda s(\lambda t(st\mathbf{F}))$ , which we denote **AND**:

$$\mathbf{AND} \equiv \lambda s(\lambda t(st\mathbf{F})).$$

One can then verify the following:

$$\begin{array}{ll} \mathbf{AND}(\mathbf{T})(\mathbf{T}) \rightarrow_{\beta} \mathbf{T} & \mathbf{AND}(\mathbf{T})(\mathbf{F}) \rightarrow_{\beta} \mathbf{F} \\ \mathbf{AND}(\mathbf{F})(\mathbf{T}) \rightarrow_{\beta} \mathbf{F} & \mathbf{AND}(\mathbf{F})(\mathbf{F}) \rightarrow_{\beta} \mathbf{F} \end{array}$$

For example,

$$\begin{aligned} \mathbf{AND}(\mathbf{T})(\mathbf{F}) &= \lambda s(\lambda t(st\mathbf{F}))(\mathbf{T})(\mathbf{F}) \\ &\rightarrow_{\beta} \lambda t(\mathbf{T}t\mathbf{F})(\mathbf{F}) \\ &\rightarrow_{\beta} \mathbf{T}\mathbf{F}\mathbf{F} \text{ (i.e., } \mathbf{T}(\mathbf{F})(\mathbf{F}).) \\ &= \lambda x(\lambda y(x))(\mathbf{F})(\mathbf{F}) \\ &\rightarrow_{\beta} \lambda y(\mathbf{F})(\mathbf{F}) \\ &\rightarrow_{\beta} \mathbf{F} \text{ (because } \mathbf{F} \text{ has no free variables, and so } \lambda y(\mathbf{F})(t) \rightarrow_{\beta} \mathbf{F} \text{ for any } t.) \end{aligned}$$

As an exercise in  $\lambda$ -calculus manipulations, you should verify the remaining three claims.

What this shows is that the  $\lambda$ -term  $\mathbf{AND}$  ‘computes’ conjunction - that is to say, given two truth values  $v_1$  and  $v_2$ , the term  $\mathbf{AND}(v_1)(v_2)$  (more typically written  $\mathbf{AND} v_1 v_2$ )  $\beta$ -reduces to the truth value of the conjunction of  $v_1$  with  $v_2$ . One can define similar  $\lambda$ -terms for other logical connectives - for example:

$$\mathbf{OR} \equiv \lambda s(\lambda t(st\mathbf{T})) \quad \mathbf{NOT} \equiv \lambda s(s\mathbf{F}\mathbf{T})$$

By combining these  $\lambda$ -terms in various ways we can of course construct more complicated connectives. What this means is that the manipulation of truth values can be interpreted as a type of manipulation (more specifically,  $\beta$ -reduction) of  $\lambda$ -terms.

The  $\lambda$ -calculus allows us to do far more than manipulate truth values. We can also encode the manipulation of natural numbers in the  $\lambda$ -calculus. Just as we began our discussion of the manipulation of truth values by encoding the truth values  $\mathbf{T}$  and  $\mathbf{F}$  with closed  $\lambda$ -terms, so too we will begin our discussion of the manipulation of natural numbers by encoding the natural numbers with closed  $\lambda$ -terms. Just as in the case of truth values, the coding is somewhat arbitrary, and there is no unique choice. The following coding of the natural numbers is most common:

$$\mathbf{0} \equiv \lambda x(\lambda y(y)), \mathbf{1} \equiv \lambda x(\lambda y(x(y))), \mathbf{2} \equiv \lambda x(\lambda y(x(x(y)))) \dots \quad (17)$$

The terms in (17) are sometimes called the *Church numerals*.

In general, we will represent the natural number  $n$  as:

$$\mathbf{n} \equiv \lambda x(\lambda y(x^n y))$$

where  $x^n y$  is  $x(x(\dots(x(y))\dots))$  with  $n$  occurrences of  $x$ . (Note the placement of parentheses.) Using this notation, you should be able to see that  $x^n(x^m y)$  is just  $x^{n+m} y$ .

With this convention, it is possible to define all sorts of complicated functions on the natural numbers. Consider, for instance, the  $\lambda$ -term  $\lambda u \lambda x \lambda y (x((u x) y))$ , which we denote **succ**. This acts on natural numbers in the following way:

$$\begin{aligned} \mathbf{succ}(\mathbf{n}) &\equiv \lambda u \lambda x \lambda y (x((u x) y))(\mathbf{n}) \\ &\rightarrow_{\beta} \lambda x \lambda y (x((\mathbf{n} x) y)) \\ &\rightarrow_{\beta} \lambda x \lambda y (x(x^n y)) \quad (\text{using } (\mathbf{n} x) y \rightarrow_{\beta} x^n y) \\ &\rightarrow_{\beta} \lambda x \lambda y (x^{n+1} y) \\ &\equiv \mathbf{n} + \mathbf{1} \end{aligned}$$

Thus we have  $\mathbf{succ} \mathbf{n} \rightarrow_{\beta} \mathbf{n} + \mathbf{1}$ , and so we may think of **succ** as a  $\lambda$ -term that ‘computes’ the successor function.

For another example, define **add** to be the  $\lambda$ -term  $\lambda u \lambda v \lambda x \lambda y [u(x)(v(x)(y))]$ . Then we have:

$$\begin{aligned} \mathbf{add}(\mathbf{n})(\mathbf{m}) &\equiv \lambda u \lambda v \lambda x \lambda y [u(x)(v(x)(y))](\mathbf{n})(\mathbf{m}) \\ &\rightarrow_{\beta} \lambda x \lambda y [\mathbf{n}(x)(\mathbf{m}(x)(y))] \\ &\rightarrow_{\beta} \lambda x \lambda y [\mathbf{n}(x)(x^m y)] \quad (\text{using } \mathbf{m}(x)(y) \rightarrow_{\beta} x^m y) \\ &\rightarrow_{\beta} \lambda x \lambda y [\lambda y (x^n y)(x^m y)] \quad (\text{using } \mathbf{n}(x) \rightarrow_{\beta} \lambda y (x^n y)) \\ &\rightarrow_{\beta} \lambda x \lambda y [x^n (x^m y)] \\ &\rightarrow_{\beta} \lambda x \lambda y [x^{n+m} y] \\ &\equiv \mathbf{n} + \mathbf{m} \end{aligned}$$

Thus, **add** represents the function of addition. It is a relatively mechanical matter to generalize these ideas to construct  $\lambda$ -terms that compute multiplication and exponentiation. You will explore this in the exercises.

In fact, more general results are possible, though they require significant effort to prove. Call a function  $f$  defined on some subset of the natural numbers *partial recursive* if it corresponds to the output of some Turing machine – that is to say, if there is some Turing machine such that for any  $n$ , if this Turing machine running on input  $n$  eventually terminates with a numerical output  $m$ , then  $f$  is defined on  $n$  and  $f(n) = m$ , while if this Turing machine running on input  $n$  does not eventually terminate with a numerical output (either because it does not terminate at all, or because it terminates with a non-numerical output on its tape), then  $f$  is undefined on  $n$ . (For a more general introduction to Turing machines and the theory of recursive functions, see chapter 3 and onwards of [3].)

For each partial recursive function  $f$  it can be shown that there is a closed  $\lambda$ -term  $\Theta_f$  such that for all natural numbers  $n$ ,  $\Theta_f(\mathbf{n})$  may be reduced to a term in  $\beta$ -normal form iff  $f(n)$  is defined. Furthermore, in the case in which  $f(n)$  is defined,

$$\Theta_f(\mathbf{n}) \twoheadrightarrow_{\beta} \mathbf{f}(\mathbf{n}).$$

(Here,  $\mathbf{f}(\mathbf{n})$  is the Church numeral corresponding to the natural number  $f(n)$ .) Thus, not only are there ways of representing the successor function, addition, multiplication and exponentiation

in the untyped  $\lambda$ -calculus  $\lambda_0$ , but also ways of representing any ‘algorithmic’ numeric function whatsoever. This result may also be straightforwardly generalized to partial recursive functions of more than one variable. For a fuller discussion and proof of this result, see chapter 4 of [4].

### Exercises for Section 1.8

1. Verify that  $\text{NOT } (\mathbf{T}) \rightarrow_{\beta} \mathbf{F}$  and  $\text{NOT } (\mathbf{F}) \rightarrow_{\beta} \mathbf{T}$ .
2. Construct a  $\lambda$ -term  $\text{ENT}$  satisfying the following axioms for entailment:
  1. For any  $\lambda$ -term  $s$ ,  $\text{ENT}(s)(\mathbf{T}) \rightarrow_{\beta} \mathbf{T}$  (i.e., everything entails  $\mathbf{T}$ ),
  2. For any  $\lambda$ -term  $s$ ,  $\text{ENT}(\mathbf{F})(s) \rightarrow_{\beta} \mathbf{T}$  (i.e., everything is entailed by  $\mathbf{F}$ ),
  3.  $\text{ENT}(\mathbf{T})(\mathbf{F}) \rightarrow_{\beta} \mathbf{F}$ .

You should verify that the term you construct has these three properties.

3. Define  $\mathbf{mult}$  to be the  $\lambda$ -term  $\lambda u \lambda v \lambda x (u(vx))$ . Verify that  $\mathbf{mult}(\mathbf{n})(\mathbf{m}) = \mathbf{p}$  where  $n \times m = p$ .
4. Construct a  $\lambda$ -term  $\mathbf{exp}$  such that  $\mathbf{exp}(\mathbf{n})(\mathbf{m}) = \mathbf{p}$  where  $n^m = p$ . Verify that the term you have constructed has this property.

## 1.9 The untyped $\lambda$ -calculus as a model of computation - II

The discussion of the previous section shows that the  $\lambda$ -calculus  $\lambda_0$  may be thought of as a model of *computation*, as powerful as any other model of computation. Just as the way in which in a Turing machine computation is implemented by a device moving along a tape, printing and deleting symbols according to definite rules, or the way in which in a modern computer computation is implemented by the manipulation of bits by logic gates, in the  $\lambda$ -calculus computation is modeled by  $\beta$ -reduction. In this way, the  $\lambda$ -calculus gives us an alternative way of thinking about the fundamental building blocks of computer programs. In particular, we can now think of all computation as ultimately reducible to the manipulation of terms in the  $\lambda$ -calculus, with the process of computation itself just becoming the process of  $\beta$ -reduction (or  $\beta\eta$ -reduction, if one wishes.) This gives a different and powerful way of thinking about the theoretical basis of computer programming.

There are texts devoted to so-called ‘programming in the  $\lambda$ -calculus’, that show how lists, more complex data structures, binary tree searches, recursion and so on can all be implemented in the  $\lambda$ -calculus. Indeed, the  $\lambda$ -calculus is perhaps the simplest example of a so-called *functional programming language*, a powerful and distinctive way of thinking about computation different from that given in the perhaps more familiar *imperative programming languages*. We will not discuss this in detail here. The interested reader should look at [6] or section 5.2 of [8] for fuller discussions.

There are several important differences between the  $\lambda$ -calculus viewed as a programming language and other more traditional programming languages. We can of course view any  $\lambda$ -term  $E$  as a program, either by (i) thinking of  $E$  itself as the starting point of a computation that terminates once we have reduced  $E$  to  $\beta$ -normal form, or (ii) by thinking of  $E$  as a (partial) function which on



input  $c$  – where  $c$  is yet another  $\lambda$ -term – produces the output corresponding to the  $\beta$ -normal form of  $E(c)$ . Much like other models of computation, in each case there is of course no guarantee that the program will halt, as not all terms have  $\beta$ -normal forms. But putting this to the side, regardless of whether we adopt point of view (i) or (ii), a curious difference with more familiar computer programs is that the outputs of these computer programs are not so much definite values, but rather  $\lambda$ -terms which are themselves further computer programs. This is different from more familiar computer programs where the output is generally some specific object, such as a number or string of characters. Still, as we have seen in the previous section, we can recover more traditional computer programs by identifying certain basic  $\lambda$ -terms with specific mathematical objects such as the truth values or the natural numbers, as in (16) or (17).

There are however further even deeper differences between the  $\lambda$ -calculus viewed as a programming language and more traditional programming languages. In a traditional computer program, a computer executes a sequence of instructions step by step. When one instruction is executed, there is no question what the next instruction to execute is - it is simply the next instruction on the list, or the next instruction as determined by some sort of rule in the presence of loops or other such programming structures. The process of  $\beta$ -reduction is however not like this. We have seen, after all, that some  $\lambda$ -terms can be  $\beta$ -reduced in different ways (though as the Church-Rosser Theorem shows, this does not mean that the program can produce different results.) So we can think of  $\lambda$ -terms as programs such that we sometimes have multiple choices as to the direction in which the program can proceed. Sometimes this is called an ‘indeterministic’ program. When we call  $\lambda$ -terms ‘programs’, they should be understood as programs in this sense.

We saw earlier that there are  $\lambda$ -terms such that some ways of  $\beta$ -reducing them lead to infinite loops, and other ways of  $\beta$ -reducing them terminate in a term in  $\beta$ -normal form. The Church-Rosser Theorem guarantees that all chains of  $\beta$ -reduction that terminate yield the same result, but that of course does not mean that all chains of  $\beta$ -reductions terminate. One might worry that in this sense the  $\beta$ -reduction of  $\lambda$ -terms is quite *unlike* traditional computing, in that we must be clever and  $\beta$ -reduce in just the right way, lest our program not terminate at all. But it turns out that this worry can be averted by following certain *protocols* when faced with the decision of how to  $\beta$ -reduce. For example, one can always perform the  $\beta$ -reduction for the redex  $\lambda x(E_1)(E_2)$  for which the  $\lambda$  is leftmost out of all redexes. It turns out that if one follows this protocol, it will terminate if there is *any* terminating sequence of  $\beta$ -reductions at all. (Demonstrating that this is so is non-trivial. See chapter 13 of [1], and in particular section 13.2 for a detailed discussion and proof.) There are other such protocols as well with this remarkable property. So long as one follows such a protocol (which can be done mechanically), the term in question reduces to a  $\beta$ -normal term if there is any such reduction at all, and one has something more like a ‘deterministic’ model of computing. A curiosity worth noting is that there are also protocols which are *unterminating* so long as there is *some* unterminating sequent of  $\beta$ -reductions!

## 1.10 Typed Computation.

We may view terms in the  $\lambda$ -calculus  $\lambda_0$  as functions that act on each other, including themselves. As we have noted already, this is somewhat different from the ordinary functions mathematicians or computer programs typically manipulate. An ordinary mathematical function like  $f(x) = 2x$  takes natural numbers or perhaps real numbers as inputs, but does not take itself as input. We do not ordinarily talk about doubling  $f$  itself, except as an indirect way to refer to something quite different, namely the function  $f(f(x)) = 4x$ . Likewise, ordinary computer programs might take as an input some natural number or string of characters, but a typical computer program does not take itself as an input.

Mathematical functions or computer programs typically act on data of a specific form and output data of a specific form. We are all familiar with functions that take natural numbers into natural numbers, or real numbers into real numbers, or vectors into real numbers, or matrices into ordered pairs of natural numbers, and so on. In each of these cases the input of the function must be of a specific form in order for the function to be applicable, and so long as the input has that specific form, the output will also have some fixed (though possibly different) form.

Given this, it is very natural to want to distinguish functions based on the structure or form of the data they input and output. A function that takes a natural number as an input and outputs a natural number is a very different sort of thing from a function that takes some sort of function as an input, and produces an ordered triple of natural numbers as an output. The untyped  $\lambda$ -calculus has no principled way of drawing these sorts of distinctions. We seek a type of formalism in which such distinctions may be drawn easily.

We begin then with the notion of a *type*. Loosely speaking, a type is some sort of domain on which a function acts, corresponding to data with some particular structure. The natural numbers  $\mathbb{N}$  is a very simple example of a type. Sometimes people refer to the type of individuals, where an individual is something which is not a function. Another very useful type is the type of truth values, which consists of the two objects **T** and **F**. This is sometimes denoted  $\mathbb{B}$ . One can also build types from other types - so we can consider the type of ordered pairs, the first element of which is an element of  $\mathbb{N}$  and the second element of which is an element of  $\mathbb{B}$ , or the type of functions from elements of  $\mathbb{N}$  to elements of  $\mathbb{B}$ . This somewhat informal understanding of the notion of a type will suffice for now.

The functions mathematicians and computer scientists work with typically not only have some sort of type as their domain, but also as their codomain. So a function from  $\mathbb{N}$  to  $\mathbb{B}$  has the type  $\mathbb{N}$  as its domain, and the type  $\mathbb{B}$  as its codomain. Consider a predicate such as ‘ $x$  is an even number’. This may be viewed as a function that takes an element of  $\mathbb{N}$  as an input, and returns an element of  $\mathbb{B}$  as an output. Thus, it may be viewed as a function from the type  $\mathbb{N}$  to the type  $\mathbb{B}$ . (For this reason, open sentences are sometimes called propositional *functions*.)

Typical computations then are *typed* – they take input of a certain type into outputs of a certain type. In addition, in typical computations quantities that appear as intermediate steps in determining the output are also typed. For example, in the course of evaluating some function  $f$  that takes  $\mathbb{N}$  to  $\mathbb{N}$  on some particular input, we might end up with some variable which always takes on

a value from  $\mathbb{B}$ , and some other variable which always takes on a value from  $\mathbb{N}$ , both of which end up being manipulated in a certain way to calculate the eventual output. The entire computation is then a manipulation of variables, each of which is typed. We call the sort of computation in which not only the input and output but also any intermediate quantities have types *typed computation*.

The untyped  $\lambda$ -calculus gives us a useful model for computation in general. Is there any similar formalism that can be used to give a model of typed computation? It turns out that there is indeed a variant of the untyped  $\lambda$ -calculus that can be used to give some sort of model of typed computation; it is unsurprisingly called the *typed*  $\lambda$ -calculus. However, perhaps unlike the case of computation in general, what counts as a typed computation and what algorithms can be captured by typed computations very much depends on how much structure one is willing to include in the class of types. In this chapter we will present the simplest version of the typed  $\lambda$ -calculus, but a sequence of richer versions of the typed  $\lambda$ -calculus will be developed in the following chapters, each version expressively richer than the previous. One might ask how well the typed  $\lambda$ -calculus and its variants serve as a model of typed computation in general, or whether there even is a most general model of typed computation in the way that Turing machines or the untyped  $\lambda$ -calculus  $\lambda_0$  give the most general possible models of general (untyped) computation. But these questions are hard and we will only be able to partially reflect on them much later.

## 1.11 Church-style Typed $\lambda$ -calculus

We begin by presenting the simplest and earliest version of the typed  $\lambda$ -calculus, in the form roughly originally developed by Church.

To start, we assume that we have some unspecified, infinite set  $\tau_1, \tau_2, \dots$  of fundamental types, out of which we will build other types. Perhaps our fundamental types consist of  $\mathbb{N}$ ,  $\mathbb{B}$ , and some other naturally occurring types - the exact details will not matter for now, and in what follows we will simply assume that  $\tau_1, \tau_2, \dots$  is a list of unspecified fundamental types. We will also assume that more complex types can be built from simpler types in a certain way. To this end, consider the definition

**Definition 1.26: Function Types.**

If  $\alpha$  and  $\beta$  are types, then  $\alpha \rightarrow \beta$  (sometimes called the ‘function type from  $\alpha$  to  $\beta$ ’) is the type of unary functions that are defined on all objects of type  $\alpha$ , and take objects of type  $\alpha$  to objects of type  $\beta$ .

So for example,  $\mathbb{N} \rightarrow \mathbb{N}$  is the type of functions that take natural numbers to natural numbers, and  $\mathbb{N} \rightarrow \mathbb{B}$  is the type of functions that take natural numbers to truth values.

We will assume that whenever we have types  $\alpha$  and  $\beta$ , we may form the type  $\alpha \rightarrow \beta$ . We will eventually add other mechanisms for building complex types out of simpler types, but the creation of function types will suffice for now. We then have the following definition

**Definition 1.27: Simple types.**

The *simple types* are those that may be obtained from the fundamental types  $\tau_1, \tau_2, \dots$  under the operation that takes two types  $\alpha$  and  $\beta$  to the type  $\alpha \rightarrow \beta$ .

For example,  $\tau_1$ ,  $(\tau_1 \rightarrow \tau_3)$  and  $(\tau_1 \rightarrow \tau_5) \rightarrow (\tau_1 \rightarrow \tau_3)$  are all simple types.

The nature of *variables* in the typed  $\lambda$ -calculus requires some discussion. In the untyped  $\lambda$ -calculus, each variable  $x, y, z, \dots$  could take *any* value. In the *typed  $\lambda$ -calculus*, by contrast, each variable will be associated with a *type*. So, for example, the variable  $x$  might represent an arbitrary element of  $\mathbb{N}$ , the variable  $y$  might represent an arbitrary element of  $\mathbb{B}$ , and so on. The easy way to implement this is to take some standard infinite list of variables  $x, y, z, \dots$  and stipulate that for each type  $\tau$ , we have the variables  $x^\tau, y^\tau, z^\tau, \dots$ , where these variables are understood as ranging over objects of type  $\tau$ . So the variables  $x^\mathbb{N}$ ,  $x^\mathbb{B}$ , and  $x^{\mathbb{N} \rightarrow \mathbb{B}}$  are all quite different things - the first is a variable that ranges over elements of  $\mathbb{N}$ , the second is a variable that ranges over elements of  $\mathbb{B}$ , and the third is a variable that ranges over elements of  $\mathbb{N} \rightarrow \mathbb{B}$ . The fact that they all involve  $x$  does not connect them in any way - they should be thought of as three completely distinct variables that have nothing more in common than  $u^\mathbb{N}, v^\mathbb{B}$  and  $w^{\mathbb{N} \rightarrow \mathbb{B}}$  do. In the typed  $\lambda$ -calculus, all variables will be typed in this way.

Just as every variable in the typed  $\lambda$ -calculus is typed, so too in fact every expression in the typed  $\lambda$ -calculus will be typed. The type of an expression will be denoted by a superscript attached to the expression itself. For example,  $[E]^\tau$  denotes the expression  $E$ , along with the stipulation that it has type  $\tau$ . Sometimes for ease of reading we will omit the brackets and simply write  $E^\tau$ . In full rigor we then define the set of terms of the typed  $\lambda$ -calculus as follows

**Definition 1.28: Simple typed  $\lambda$ -terms.**

The set of *simply typed  $\lambda$ -terms* are defined inductively as follows:

- (i) For any variable  $v$  and any simple type  $\tau$ ,  $v^\tau$  is a simply typed  $\lambda$ -term of type  $\tau$ . (To reduce clutter, we will generally not write this as  $[v]^\tau$ ).
- (ii) For simple types  $\alpha$  and  $\beta$ , if  $[M]^{\alpha \rightarrow \beta}$  and  $[N]^\alpha$  are simply typed  $\lambda$ -terms, then  $[[M]^{\alpha \rightarrow \beta}([N]^\alpha)]^\beta$  is a simply typed  $\lambda$ -term.
- (iii) For simple types  $\alpha$  and  $\beta$ , if  $[M]^\beta$  is a simply typed  $\lambda$ -term and  $x^\alpha$  a variable, then  $[\lambda x^\alpha([M]^\beta)]^{\alpha \rightarrow \beta}$  is a simply typed  $\lambda$ -term.

Terms of the typed  $\lambda$ -calculus then look very much like terms of the untyped  $\lambda$ -calculus, with superscripts added to each subexpression. Examples of terms in the typed  $\lambda$ -calculus include:

$$\begin{aligned}
 & [x^{\alpha \rightarrow \beta}(y^\alpha)]^\beta \\
 & [[\lambda x^\alpha(x^\beta)]^{\alpha \rightarrow \beta}(y^\alpha)]^\beta \\
 & z^{(\alpha \rightarrow \beta) \rightarrow \gamma}([\lambda x^\alpha(x^\beta)]^{\alpha \rightarrow \beta})
 \end{aligned} \tag{18}$$

The clauses of Definition 1.28 might look complex, but actually express fairly simple ideas. Clause (i) simply gives us an inexhaustible list of variables of each type. Clause (ii) corresponds to the idea of *application*, this time in the context of the typed  $\lambda$ -calculus. In the untyped lambda calculus, we could apply any function to any input - that is to say,  $M(N)$  was always a valid  $\lambda$ -term for any  $\lambda$ -terms  $M$  and  $N$ . That however is no longer the case. To apply  $M$  to  $N$ ,  $M$  must be a function and  $N$  must be an object of the right type to act as an input to  $M$ . More specifically, for some types  $\alpha$  and  $\beta$ ,  $M$  must be a term of type  $\alpha \rightarrow \beta$  and  $N$  must be a term of type  $\alpha$ . In such a case, the result of applying  $M$  to  $N$  will have type  $\beta$ . Such a requirement is the idea behind clause (ii). So for example, if  $\alpha$  and  $\beta$  are distinct simple types, the term  $x^{\alpha \rightarrow \beta}(y^\beta)$  is *not* a legal simply typed  $\lambda$ -term, as it involves trying to apply a function from  $\alpha$  to  $\beta$  to an input from  $\beta$ . The term  $[x^{\alpha \rightarrow \beta}(y^\alpha)]^\beta$  is by contrast acceptable and has overall type  $\beta$  because it is the result of applying a function of type  $\alpha \rightarrow \beta$  to an input of type  $\alpha$ , which of course yields an output of type  $\beta$ .

Consider now clause (iii). In the untyped context, the term  $\lambda x(M)$  corresponded to treating the expression  $M$  as a function of some variable  $x$ . This is a function which on input  $c$  yields  $M[x/c]$ . Suppose now that  $x$  has some type  $\alpha$  and  $M$  has some type  $\beta$ . Then we can think of  $M[x/c]$  as a function that takes an input of type  $\alpha$  and produces an output of type  $\beta$ . Thus, in the typed context the  $\lambda$ -term  $\lambda x(M)$  - written more fully as  $\lambda x^\alpha([M]^\beta)$  - has the simple type  $\alpha \rightarrow \beta$ . That is the idea behind clause (iii).

We will often omit not only the brackets  $[\cdot]$ , but also type symbols when they are obvious or can be easily inferred from the context. So for instance, rather than writing  $[[M]^{\alpha \rightarrow \beta}([N]^\alpha)]^\beta$ , we will typically just write  $M^{\alpha \rightarrow \beta}(N^\alpha)$  (or even just  $M^{\alpha \rightarrow \beta}N^\alpha$ ), as the fact that this term has type  $\beta$  follows trivially from the fact that it involves applying a function of type  $\alpha \rightarrow \beta$  to an object of type  $\alpha$ . Likewise, rather than writing  $[\lambda x^\alpha(M^\beta)]^{\alpha \rightarrow \beta}$ , we will typically write  $\lambda x^\alpha(M^\beta)$  (or even just  $\lambda x^\alpha M^\beta$ .) The terms (18) can therefore be written more simply:

$$\begin{aligned} & x^{\alpha \rightarrow \beta}(y^\alpha) \\ & \lambda x^\alpha(x^\beta)(y^\alpha) \\ & z^{(\alpha \rightarrow \beta) \rightarrow \gamma}(\lambda x^\alpha(x^\beta)) \end{aligned}$$

In general, when a  $\lambda$ -term is written out in full it suffices to indicate only the type of each variable, as the type of any larger expression can then be easily deduced, as seen in our three examples.

In addition, for each bound variable it typically suffices to indicate the type of the variable only when it is introduced with the corresponding  $\lambda$ -operator. For example, the term  $\lambda x^\alpha \lambda y^\beta(x^\alpha y^\beta)$  can be rewritten  $\lambda x^\alpha \lambda y^\beta(xy)$ , as the types of  $x$  and  $y$  have already been stipulated to be  $\alpha$  and  $\beta$  respectively. There are, however, situations where this convention can lead to confusion. For example, if we were to start with the typed  $\lambda$ -term  $\lambda x^\alpha \lambda x^\beta(x^\alpha x^\beta)$  and use the convention in question, we would end up with  $\lambda x^\alpha \lambda x^\beta(xx)$ , which is potentially confusing (what are the types of the unadorned  $x$ s?) This sort of situation can be avoided by renaming variables and working with  $\alpha$ -equivalent terms. So  $\lambda x^\alpha \lambda x^\beta(x^\alpha x^\beta)$  is  $\alpha$ -equivalent to  $\lambda x^\alpha \lambda y^\beta(x^\alpha y^\beta)$ , which may be abbreviated  $\lambda x^\alpha \lambda y^\beta(xy)$ , and all ambiguity has been removed. In general, we will not omit information when it could potentially lead to ambiguity.

With all these conventions in mind, the following may be regarded as fully typed terms:

$$\begin{aligned} &\lambda x^{\alpha \rightarrow \beta}(x(z^\alpha)) \\ &\lambda x^\alpha \lambda y^{\alpha \rightarrow \beta}(yx) \\ &\lambda x^{\alpha \rightarrow \beta} \lambda z^{(\alpha \rightarrow \beta) \rightarrow \gamma}(z(\lambda y^\alpha(xy))) \end{aligned}$$

In particular, given the type of each free and bound variable in a  $\lambda$ -term, the type of every subexpression (including the type of the overall term itself) can be determined.

In fact, sometimes even less information is necessary to determine the type of a  $\lambda$ -term and all its subexpressions. For instance, consider the term

$$\lambda x^{\alpha \rightarrow \beta}(xy).$$

Here the bound variable  $x$  is typed, but the open variable  $y$  has not been explicitly typed. However, from the fact that there is a subexpression of the form  $xy$  and  $x$  has type  $\alpha \rightarrow \beta$ , it follows that the type of  $y$  must be  $\alpha$ . Thus the expression  $\lambda x^{\alpha \rightarrow \beta}(xy)$  can be regarded as fully typed, as it contains enough information for us to deduce the type of every subexpression. Nevertheless, in order to make things as simple as possible for the reader, we will generally explicitly indicate the types of all bound and free variables in a typed term, writing for example

$$\lambda x^{\alpha \rightarrow \beta}(xy^\alpha).$$

As in the untyped  $\lambda$ -calculus, the  $\lambda$ -term  $\lambda x^\sigma(M^\tau)(N^\sigma)$  simplifies to

$$M^\tau[N^\sigma/x^\sigma].$$

As before, we will call the act of simplifying  $\lambda x^\sigma(M^\tau)(N^\sigma)$  to  $M^\tau[N^\sigma/x^\sigma]$   $\beta$ -reduction. We use the same notational conventions as before:

**Definition 1.29:  $\beta$ -reduction.**

An expression of the form  $\lambda x^\sigma(M^\tau)(N^\sigma)$  (where  $x^\sigma$  is a variable and  $M^\tau, N^\sigma$  are  $\lambda$ -terms) is called a  $\beta$ -redex (or just a redex). We have

- (i)  $A \rightarrow_\beta B$  iff  $A$  has the form  $\tau_1 \lambda x^\sigma(M^\tau)(N^\sigma) \tau_2$  and  $B$  has the form  $\tau_1 M^\tau[N^\sigma/x^\sigma] \tau_2$ .
- (ii)  $A \rightarrow_\beta B$  ( $A$   $\beta$ -reduces to  $B$ ) iff for some sequence  $X_0, X_1, \dots, X_n$  of  $\lambda$ -terms we have that  $A = X_0, B = X_n$ , and  $X_0 \rightarrow_\beta X_1 \dots \rightarrow_\beta X_n$ . This includes the degenerate case in which  $A$  and  $B$  are identical.
- (iii) A term  $A$  is in  $\beta$ -normal form just in case no  $\beta$ -reduction may be performed on it.

Analogous terminology is used for  $\eta$ -reduction and  $\beta\eta$ -reduction.

We use the term  $\lambda_0^{\text{type}}$  to refer to this version of the  $\lambda$ -calculus, with the set of rules of formation of typed  $\lambda$ -terms and rules of reduction just presented. Sometimes it is also called a *Church style typed  $\lambda$ -calculus*. We present its rules together for convenience:

**Definition 1.30: Rules of the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ .**Types

The *simple types* are those that may be obtained from the fundamental types  $\tau_1, \tau_2, \dots$  under the operation that takes two types  $\alpha$  and  $\beta$  to the type  $\alpha \rightarrow \beta$ .

Terms

The set of *simply typed  $\lambda$  terms* are defined inductively as follows:

- (i) For any variable  $v$  and any simple type  $\tau$ ,  $v^\tau$  (or  $[v]^\tau$ ) is a simply typed  $\lambda$ -term.
- (ii) For simple types  $\alpha$  and  $\beta$ , if  $[M]^{\alpha \rightarrow \beta}$  and  $[N]^\alpha$  are simply typed  $\lambda$ -terms, then  $M^{\alpha \rightarrow \beta}(N^\alpha)$  (or more fully,  $[[M]^{\alpha \rightarrow \beta}([N]^\alpha)]^\beta$ ) is a simply typed  $\lambda$ -term.
- (iii) For simple types  $\alpha$  and  $\beta$ , if  $[M]^\beta$  is a simply typed  $\lambda$ -term and  $x^\alpha$  a variable, then  $\lambda x^\alpha(M^\beta)$  (or more fully,  $[\lambda x^\alpha([M]^\beta)]^{\alpha \rightarrow \beta}$ ) is a simply typed  $\lambda$ -term.

Reduction Rules

- (i)  $\lambda x^\sigma([M]^\tau)([N]^\sigma) \rightarrow_\beta [M]^\tau[[N]^\sigma/x^\sigma]$
- (ii)  $\lambda x^\sigma([E]^\tau(x^\sigma)) \rightarrow_\eta [E]^\tau$ , where  $x^\sigma$  does not appear free in  $[E]^\tau$ .

The typed  $\lambda$ -calculus is quite different from the untyped  $\lambda$ -calculus in several important ways. Unlike the untyped  $\lambda$ -calculus  $\lambda_0$ , it may be shown that every term of the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  can be  $\beta$ -reduced to a term in  $\beta$ -normal form (and moreover can be  $\beta\eta$ -reduced to a term in  $\beta\eta$ -normal form.) Recall the problematic *untyped*  $\lambda$ -term:

$$\lambda x(x(x))(\lambda y(y(y))).$$

This *untyped* term has no  $\beta$ -normal form - when we  $\beta$ -reduce it, we simply get the original term back. However, there is no term like this in the *typed*  $\lambda$ -calculus. The sub-term  $x(x)$  of our problematic untyped  $\lambda$ -term has nothing corresponding to it in the typed  $\lambda$ -calculus. In particular, for any type  $\alpha$ ,  $x^\alpha(x^\alpha)$  is not a legal term, as it does not conform to the pattern in clause (ii) of the definition of  $\lambda$ -terms in Definition 1.30. In order for  $x^\beta(x^\alpha)$  to be a legal term in the typed  $\lambda$ -calculus,  $\beta$  must have the form  $\alpha \rightarrow \gamma$  for some  $\gamma$ . In order for  $x^\alpha(x^\alpha)$  to be a legal term,  $\alpha$  itself would then have to have the form  $\alpha \rightarrow \gamma$ , and no simple type (or any type we consider in this book) has this sort of self-referential character.

In the exercises, you will be walked through a proof that each term in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  has a  $\beta$ -normal form. Actually, it turns out that something stronger is true – for each term in  $\lambda_0^{\text{type}}$  every series of  $\beta$ -reductions eventually terminates with a term in  $\beta$ -normal form. (This property is sometimes called *strong normalization*.) That is to say, not only is it the case that for every typed  $\lambda$ -term of  $\lambda_0^{\text{type}}$  there is *some* sequence of  $\beta$ -reductions taking it to a term in  $\beta$ -normal form, but in fact *every* sequence of  $\beta$ -reductions eventually leads to a term in  $\beta$ -normal form. That is, we have the following:

**Theorem 1.31: Strong Normalizability of  $\lambda_0^{\text{type}}$  under  $\beta$ -reduction.**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$$

The proof of this theorem is extremely delicate and is included in the appendix to this chapter. As before, we also have a version of the Church-Rosser theorem for the typed  $\lambda$ -calculus:

**Theorem 1.32: The Church-Rosser Theorem for  $\beta$ -reductions in  $\lambda_0^{\text{type}}$ .**

In the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , if  $X \rightarrow_{\beta} Y_1$  and  $X \rightarrow_{\beta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ .

We also discuss this in the appendix. Analogs of these theorems also hold for  $\beta\eta$ -reductions in  $\lambda_0^{\text{type}}$ . In particular, we have the following

**Theorem 1.33: Strong Normalizability of  $\lambda_0^{\text{type}}$  under  $\beta\eta$ -reduction.**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta\eta} \tau_2 \rightarrow_{\beta\eta} \tau_3 \rightarrow_{\beta\eta} \dots$$

(where  $X \rightarrow_{\beta\eta} Y$  iff  $X \rightarrow_{\beta} Y$  or  $X \rightarrow_{\eta} Y$ .)

**Theorem 1.34: The Church-Rosser Theorem for  $\beta\eta$ -reductions in  $\lambda_0^{\text{type}}$ .**

In the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

Theorem 1.33 actually follows easily from Theorem 1.31, using the fact proved in an earlier exercise that if a term has an infinite  $\beta\eta$ -reduction sequence, then it also has an infinite  $\beta$ -reduction sequence. (Although the argument given for this was in the context of the untyped  $\lambda$ -calculus, the addition of types does not affect the argument.) The proof of Theorem 1.34 requires more attention, and is left to the appendix.

As before, it follows easily from the Church-Rosser Theorem(s) for  $\lambda_0^{\text{type}}$  that starting with the same  $\lambda_0^{\text{type}}$  term, different sequences of  $\beta$ -reductions (or  $\beta\eta$ -reductions) that terminate with



a term in  $\beta$ -normal form (or  $\beta\eta$ -normal form) must in fact terminate with the same term (up to  $\alpha$ -equivalence.)

Putting all these results together we have that starting with a fixed term of the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , every sequence of  $\beta$ -reductions (or  $\beta\eta$ -reductions) eventually terminates, and terminates with exactly the same  $\beta$ -normal form (or  $\beta\eta$ -normal form) term.

The fact that every sequence of  $\beta$ -reductions eventually terminates makes the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  a very convenient programming language, as we do not have to worry about programs falling into infinite loops, for example. But this comes at a cost. We cannot hope to simulate all Turing Machines in the typed  $\lambda$ -calculus, as not all Turing Machines halt. If one thinks that there is some intuitive sense in which the running of a Turing machine program is an example of typed computation, then one will have to concede that there are examples of typed computation that are *not* captured by the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  (or even any of the extensions we will consider in the next chapters.) Perhaps then the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  should be taken a model of something more like *convergent* (i.e., terminating) typed calculation. While there is something right about this, the typed  $\lambda$ -calculus does not seem to be a ‘universal’ model of even convergent typed calculation in the way the untyped  $\lambda$ -calculus is something like a universal model of general computation, as by adding further structure to the types of the untyped  $\lambda$ -calculus (as we will in the following chapters) we get increasing expressive power. This is unlike the case of the untyped  $\lambda$ -calculus in which we can already represent all computable functions, and so adding further computational machinery does not end up increasing expressive power. Indeed, it is not clear at present whether there is any universal model of typed computation or even typed convergent computation analogous to the universal model of general computation given by Church’s Thesis.

The important point for now is just that the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , in contrast with the untyped  $\lambda$ -calculus  $\lambda_0$ , is a somewhat limited programming language. The Curry-Howard Correspondence will reveal even further limitations. Still, the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  is not *hopelessly* limited, as there are many aspects of data manipulation that can be captured by it. Moreover, in further chapters we will add various resources to the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  to significantly expand its expressive power, while at the same time keeping everything fully ‘typed’. So both in its own right and as a starting point for further expansion, the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  is both a programming language and mathematical object of great interest.

### Exercises for Section 1.11

1. Say whether the following terms are grammatical or ungrammatical terms of the Church style typed  $\lambda$ -calculus. For each grammatical term, if it is not already in  $\beta$ -normal form you should (i) reduce it to  $\beta$ -normal form, and (ii) state what the type of the  $\beta$ -normal form term is. You should assume that  $A$  and  $B$  are distinct types.

(i).  $x^A(y^B)$

(ii).  $\lambda x^A(y^B)(z^B)$

- (iii).  $\lambda x^A(y^B)(z^A)$
- (iv).  $\lambda x^A(x^A)(\lambda y^B(y^B))$
- (v).  $\lambda x^A(x^A)(\lambda y^A(y^A))$
- (vi).  $\lambda x^{A \rightarrow A}(x^{A \rightarrow A})(\lambda y^A(y^A))$
- (vii).  $\lambda x^{A \rightarrow A}(x^{A \rightarrow A})(\lambda y^A(y^A))(z^A)$

2. (a) Prove that if the type of each free and bound variable in a  $\lambda$ -term  $M$  is given, that the type of every subexpression of  $M$  can be determined. (Hint: proceed by induction on the construction of  $\lambda$ -terms.)

(b) Is it true that if the type of a  $\lambda$ -term  $M$  as well as the type of each free variable is given, that the type of every subexpression of  $M$  can be determined? If so, prove this fact. If not, give a counterexample.

(c) Is it true that if the type of a  $\lambda$ -term  $M$  as well as the type of each bound variable is given, that the type of every subexpression of  $M$  can be determined? If so, prove this fact. If not, give a counterexample.

3. We prove in this exercise that there is a strategy for  $\beta$ -reducing terms in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  that always terminates. Thus, every term in the typed  $\lambda$ -calculus has a normal form. (This is a weaker result than Theorem 1.31, but is easier to prove.)

First of all, we define the *complexity*  $c(\tau)$  of a simple type inductively by stipulating that  $c(\tau)$  is 0 if  $\tau$  is one of the fundamental types, and  $c(\tau)$  is  $1 + \max\{c(\sigma_1), c(\sigma_2)\}$  if  $\tau$  is a type of the form  $\sigma_1 \rightarrow \sigma_2$ .

We next associate with each  $\lambda$ -term  $s$  a pair of natural numbers  $n(s)$  and  $m(s)$ , where  $n(s)$  is the maximum value of  $c(\tau)$  for any variable  $x^\tau$  that appears bound in  $s$ , and  $m(s)$  is the number of occurrences in  $s$  of sub-expressions of the form  $\lambda x^\tau(\dots)$ , where  $\tau$  is a type with  $c(\tau) = n(s)$ . So loosely speaking,  $n(s)$  is the highest complexity of the type of any bound variable in  $s$ , and  $m(s)$  is the number of abstractions in  $s$  over a variable whose type has that maximum complexity.

Argue that if  $s$  is a typed  $\lambda$ -term which is not in  $\beta$ -normal form, then there is a term  $s'$  that may be obtained from  $s$  by a single  $\beta$ -reduction (i.e.,  $s \rightarrow_\beta s'$ ) such that either (i)  $n(s') < n(s)$ , or (ii)  $n(s') = n(s)$  and  $m(s') < m(s)$ . Argue from this that for any typed  $\lambda$ -term  $s$ , there is *some* sequence of  $\beta$ -reductions beginning with  $s$  that terminates. You should in fact be able to describe an algorithm for the  $\beta$ -reduction of typed  $\lambda$ -calculus terms that always terminates.

4. Can the strategy of the previous problem be adapted to give a strategy for  $\beta\eta$ -reducing terms in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  that always terminates?

## 1.12 Curry-style Typed $\lambda$ -calculus

Somewhat confusingly, there are several distinct formalisms that go under the name of the typed  $\lambda$ -calculus. First, there is the *Church-style* typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  that we have already seen. Second, there is the so-called *Curry-style* typed  $\lambda$ -calculus, of which several distinct versions exist. We

describe one version of the Curry-style typed  $\lambda$ -calculus this section and another in the next.

In the Church-style  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , each  $\lambda$ -term has a single, unique type. However, in the first version of the Curry style  $\lambda$ -calculus that we consider, the terms are ordinary terms of the *untyped*  $\lambda$ -calculus  $\lambda_0$ , each of which is assigned a *set* of possible types. To see why we might want to do this, consider the untyped  $\lambda$ -term  $\lambda x(x)$ . This of course is the identity function. It is built up from the variable  $x$ . If we assign  $x$  the type  $\mathbb{N}$ , then the untyped  $\lambda$ -term  $\lambda x(x)$  is transformed into the typed  $\lambda$ -term  $\lambda x^{\mathbb{N}}(x^{\mathbb{N}})$  (or  $\lambda x^{\mathbb{N}}(x)$ , using our abbreviations). This of course is just the identity function that takes an element of  $\mathbb{N}$  to itself, and has type  $\mathbb{N} \rightarrow \mathbb{N}$ . It is a term of the Church-style typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ . Likewise, if we think of  $x$  as having the type  $\mathbb{B}$ , then the untyped  $\lambda$ -term  $\lambda x(x)$  is transformed into the typed  $\lambda$ -term  $\lambda x^{\mathbb{B}}(x^{\mathbb{B}})$  (or  $\lambda x^{\mathbb{B}}(x)$ ), which is the identity function that takes an element of  $\mathbb{B}$  to itself, and has type  $\mathbb{B} \rightarrow \mathbb{B}$ . In general, for any type  $\tau$ , by thinking of  $x$  as having type  $\tau$  our original term  $\lambda x(x)$  can be ‘typed’ into the term  $\lambda x^{\tau}(x^{\tau})$  (or  $\lambda x^{\tau}(x)$ ) of the Church-style typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , where it acquires the type  $\tau \rightarrow \tau$ .

Because it is possible to add superscripts to the term  $\lambda x(x)$  in different ways to produce different terms of the Church-style typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , we can think of the untyped term  $\lambda x(x)$  as having many possible types. The set of types it is possible for the term  $\lambda x(x)$  to acquire in this way is of course just the set of types of the form  $X \rightarrow X$ .

Just as we can examine the set of possible types a closed  $\lambda$ -term like  $\lambda x(x)$  can be given, so too we can examine the set of possible types an open  $\lambda$ -term can be given. Consider for example the open untyped  $\lambda$ -term  $\lambda x(y)$ . By assigning  $x$  the type  $\sigma$  and  $y$  the type  $\tau$ , the term  $\lambda x(y)$  is transformed into the term  $\lambda x^{\sigma}(y^{\tau})$ , a term of  $\lambda_0^{\text{type}}$ , which of course has the type  $\sigma \rightarrow \tau$ . In this way, the set of types it is possible for the term  $\lambda x(y)$  to acquire is the set of types of the form  $X \rightarrow Y$ .

In determining the possible types of an expression, we are constrained by the fact that when a typed expression  $M^{\sigma}$  is applied to a typed expression  $N^{\tau}$ , the resulting expression  $M^{\sigma}(N^{\tau})$  makes grammatical sense only if  $\sigma$  is of the form  $\tau \rightarrow \gamma$ . As a result of this constraint, the set of possible types a given untyped term can assume can be empty. For example, if  $x$  is a variable, consider the open term  $x(x)$ . If we assign  $x$  the type  $\sigma$ , it has the form  $x^{\sigma}(x^{\sigma})$ . But this is ungrammatical; if the second  $x$  has type  $\sigma$ , then the first  $x$  needs to have type  $\sigma \rightarrow \tau$  for some  $\tau$  in order for the whole expression to make sense. So it is impossible to assign  $x(x)$  any type at all in a consistent manner, and thus its set of possible types is empty. Any untyped term in which  $x(x)$  occurs (such as the closed term  $\lambda x(x(x))$ ) will also have an empty set of possible types for the same reason.

Let us consider another example. Take the  $\lambda$ -term  $\lambda x(\lambda y(x))$ . What types can this be assigned? Suppose  $x$  has type  $\tau$ , and  $y$  has type  $\sigma$ . Then  $\lambda y(x)$  is a function from objects of type  $\sigma$  to objects of type  $\tau$ , and so is itself a term of type  $\sigma \rightarrow \tau$ . The entire expression  $\lambda x(\lambda y(x))$  is thus a function which takes an object of type  $\tau$  into an object of type  $\sigma \rightarrow \tau$ , and so has type  $\tau \rightarrow (\sigma \rightarrow \tau)$ . We have thus shown that the untyped term  $\lambda x(\lambda y(x))$  can be ‘typed’ into the following term of  $\lambda_0^{\text{type}}$

$$[\lambda x^{\tau}(\lambda y^{\sigma}(x^{\tau}))]^{\sigma \rightarrow \tau}{}^{\tau \rightarrow (\sigma \rightarrow \tau)}$$

or more simply,

$$\lambda x^{\tau}(\lambda y^{\sigma}(x)).$$

Although this might seem pedantic, we would like to formalize this sort of reasoning about which types an untyped  $\lambda$ -term can be assigned. Central to this sort of reasoning will be claims of the following form, for an untyped  $\lambda$ -term  $E$

Consistent with the constraint that the variables  $v_1, v_2, \dots, v_n$  (which include all the free variables of  $E$ ) be assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, it is possible to type  $E$  so that it becomes a  $\lambda_0^{\text{type}}$  term of type  $\sigma$ .

We express this using the following notation

**Definition 1.35: Type Theoretic Sequents.**

For an untyped  $\lambda$ -term  $E$  and list of variables  $v_1, v_2, \dots, v_n$  that include (but is not necessarily limited to) all the free variables of  $E$ , we write

$$v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n \vdash E : \sigma$$

to mean that consistent with the constraint that the free variables  $v_1, v_2, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, it is possible to type  $E$  so that it becomes a  $\lambda_0^{\text{type}}$  term of type  $\sigma$ .

The claim

$$v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n \vdash E : \sigma$$

is called a *type theoretic sequent* (or just a *sequent*), and its left hand side

$$v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$$

is its *context* or *environment*.

A sequent then just says that in a given context, a term may be assigned a given type. The sequent  $\vdash E : \sigma$  with empty context then just says that  $E$  is a closed term that can be assigned the type  $\sigma$  – i.e.,  $E$  has  $\sigma$  among its possible types.

Using this notation, the examples we have give so far in this section may be rewritten

$$\begin{aligned} &\vdash \lambda x(x) : \tau \rightarrow \tau \\ y : \tau &\vdash \lambda x(y) : \sigma \rightarrow \tau \\ y : \sigma &\vdash \lambda x(\lambda y(x)) : \tau \rightarrow (\sigma \rightarrow \tau) \end{aligned}$$

for any types  $\sigma$  and  $\tau$ .

Within a context  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$  we allow a variable to appear more than once. If it appears more than once, however, it must be assigned the same type each time. So for example,

$$x : \tau, x : \tau, y : \sigma$$

is a grammatical context, but

$$x : \tau, x : \sigma, y : \sigma$$

is not a grammatical context if  $\sigma$  and  $\tau$  are distinct.

The context can be thought of as consisting of a sequence of ‘variable declarations’. For example, we can think of the context:

$$x : \mathbb{N}, y : \mathbb{N} \rightarrow \mathbb{N}$$

as expressing:

$x$  is a natural number, and  $y$  is a function from natural numbers to natural numbers.

or, from a computer programmer’s point of view

let  $x$  be a natural number, and let  $y$  be a function from natural numbers to natural numbers.

With all this understood, we can then easily write down some ‘laws’ that sequents of the form  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n \vdash E : \sigma$  must obey. In particular, we have:

**Definition 1.36: Rules of the typed  $\lambda$ -calculus  $\mathbf{TR}_0$ .**

- (i)  $\Gamma, v : \tau \vdash v : \tau$  (Var)
- (ii)  $\frac{\Gamma \vdash E_1 : \tau \rightarrow \sigma \quad \Gamma \vdash E_2 : \tau}{\Gamma \vdash E_1(E_2) : \sigma}$  (App)
- (iii)  $\frac{\Gamma, v : \tau \vdash E : \sigma}{\Gamma \vdash \lambda v(E) : \tau \rightarrow \sigma}$  (Abs)

where  $v$  is a variable,  $E_1, E_2$  and  $E$  are terms of the untyped  $\lambda$ -calculus,  $\sigma, \tau$  are types, and all contexts are assumed to be grammatical.

‘Var’ is here an abbreviation for *variable rule*, ‘App’ for *application rule*, and ‘Abs’ for *abstraction rule*. Sometimes these rules are called *typing rules*, and with this in mind we call the system of three rules just given  $\mathbf{TR}_0$ . (Other texts sometimes use the name  $\lambda_{\rightarrow}$  for this set of rules.) Our first version of the Curry-style typed  $\lambda$ -calculus is constituted by this set of typing rules.

The (Var) rule states the obvious fact that in some context in which the variable  $v$  has been assigned type  $\tau$ ,  $\tau$  is a possible type of  $v$ . The (App) rule states that if  $E_1$  has possible type  $\tau \rightarrow \sigma$  in some context  $\Gamma$  and  $E_2$  has possible type  $\tau$  in the same context  $\Gamma$ , then in that context  $E_1(E_2)$  has possible type  $\sigma$ . The (Abs) rule states that if in some context  $\Gamma, v : \tau$  the term  $E$  has the possible type  $\sigma$ , then in the context  $\Gamma$ , the term  $\lambda v(E)$  has possible type  $\tau \rightarrow \sigma$ . These three rules may thus be read as statements about the structure of the set of possible types that each term can be assigned.

We prove that an untyped  $\lambda$ -term can be assigned a certain type by constructing a ‘proof’ of this fact using the above rules. More specifically, to show that  $E$  can be assigned the type  $X$  it suffices to construct a proof of

$$\Gamma \vdash E : X$$

in  $\mathbf{TR}_0$  for some grammatical context  $\Gamma$ .

To take our most basic example, to show that the  $\lambda$ -term  $\lambda x(x)$  can be assigned the type  $\tau \rightarrow \tau$  for any type  $\tau$ , we have the following argument:

$$\frac{x : \tau \vdash x : \tau \quad (\text{Var})}{\vdash \lambda x(x) : \tau \rightarrow \tau} (\text{Abs})$$

Note that while this argument shows that every type of the form  $\tau \rightarrow \tau$  is a possible type of  $\lambda x(x)$ , it does not show that every possible type of  $\lambda x(x)$  has the form  $\tau \rightarrow \tau$ . While true, this latter fact requires a different sort of argument that we will not consider here.

For another example, to show that the expression  $\lambda x(\lambda y(x))$  can be assigned the type  $\tau \rightarrow (\sigma \rightarrow \tau)$  for any types  $\tau$  and  $\sigma$ , we have the following argument:

$$\frac{\frac{x : \tau, y : \sigma \vdash x : \tau \quad (\text{Var})}{x : \tau \vdash \lambda y(x) : \sigma \rightarrow \tau} (\text{Abs})}{\vdash \lambda x(\lambda y(x)) : \tau \rightarrow (\sigma \rightarrow \tau)} (\text{Abs})$$

This argument shows that every type of the form  $\tau \rightarrow (\sigma \rightarrow \tau)$  is a possible type of  $\lambda x(\lambda y(x))$ .

For a more complex example, to show that the expression  $\lambda x(\lambda y(x(xy)))$  can be assigned any type of the form  $(\sigma \rightarrow \sigma) \rightarrow (\sigma \rightarrow \sigma)$ , we have the following derivation:

$$\frac{\frac{\frac{x : \sigma \rightarrow \sigma, y : \sigma \vdash x : \sigma \rightarrow \sigma \quad x : \sigma \rightarrow \sigma, y : \sigma \vdash y : \sigma}{x : \sigma \rightarrow \sigma, y : \sigma \vdash xy : \sigma} (\text{App})}{x : \sigma \rightarrow \sigma, y : \sigma \vdash x(xy) : \sigma} (\text{Abs})}{x : \sigma \rightarrow \sigma \vdash \lambda y(x(xy)) : \sigma \rightarrow \sigma} (\text{Abs})}{\vdash \lambda x(\lambda y(x(xy))) : (\sigma \rightarrow \sigma) \rightarrow (\sigma \rightarrow \sigma)} (\text{Abs})$$

Each ‘leaf’ of this tree is of the form (Var). We call arguments like these *typing derivations*.

In typing derivations, we allow ourselves to ‘silently’ replace untyped  $\lambda$ -terms by  $\alpha$ -equivalent terms, without requiring any explicit step in which the relevant bound variables are renamed. So for example, consider the untyped  $\lambda$ -term

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

As discussed earlier, we regard the innermost  $(x)$  as bound by the second  $\lambda x$  from the left rather than the first  $\lambda x$  from the left. (In general, when there are multiple options for binding a variable, the variable is bound by the closest possible binding operator that is grammatically capable of binding it.) Given this, we should expect that for any types  $\sigma$  and  $\tau$  this  $\lambda$ -term can be assigned the type  $\sigma \rightarrow (\tau \rightarrow \tau)$ . Proving this is easy:

$$\frac{\frac{x : \sigma, y : \tau \vdash y : \tau \quad (\text{Var})}{x : \sigma \vdash \lambda y(y) : \tau \rightarrow \tau} (\text{Abs})}{\vdash \lambda x(\lambda x(x)) : \sigma \rightarrow (\tau \rightarrow \tau)} (\text{Abs})$$

where in the last step the untyped  $\lambda$ -term  $\lambda x(\lambda y(y))$  has been silently replaced by the  $\alpha$ -equivalent term  $\lambda x(\lambda x(x))$ . You should be able to convince yourself that without silent renaming of this sort, no proof of the sequent  $\vdash \lambda x(\lambda x(x)) : \sigma \rightarrow (\tau \rightarrow \tau)$  would be possible.

To avoid having to rename variables in the middle of a proof, it can be useful to start with  $\lambda$ -terms in which all bound variables have different names, and no variable is both free and bound. Sometimes this is called the *Barendregt variable convention*.

**Definition 1.37: Barendregt Variable Convention.**

Say that a  $\lambda$ -term  $E$  obeys the *Barendregt variable convention* iff all bound variables of  $E$  have different names, and no variable of  $E$  occurs both free and bound in  $E$ .

Note that sometimes this convention is defined to merely require that no variable of  $E$  occurs both free and bound in  $E$ , but leaves open the possibility that distinct bound variables have the same name, as in the terms  $\lambda x(\lambda x(y))$  or  $\lambda x(x)(\lambda x(x))$ . We formulate the convention in a more strict way to exclude these possibilities.

Note also that in the Barendregt variable convention, the  $\lambda$ -term  $E$  is not being thought of as an equivalence class of  $\lambda$ -terms, but rather as a particular element from this class, i.e., a specific expression. This is an example of our casually moving between talk of  $\lambda$ -terms as specific strings of symbols (syntactic objects) and talk of  $\lambda$ -terms as equivalence classes of such syntactic objects.

The proof of the following theorem is straightforward and is left to the exercises.

**Theorem 1.38**

If  $E$  is a  $\lambda$ -term  $E$  obeying the Barendregt variable convention and  $E'$  is a subterm of  $E$ , then  $E'$  obeys the Barendregt variable convention.

Given any  $\lambda$ -term, there is always an  $\alpha$ -equivalent  $\lambda$ -term that obeys the Barendregt variable convention. For instance, in  $x(\lambda x(x))$  there are both free and bound occurrences of  $x$ , and in  $\lambda x(\lambda x(zx))$  we have two declarations of the bound variable  $x$ . However,  $x(\lambda x(x))$  is  $\alpha$ -equivalent to  $x(\lambda y(y))$ , and  $\lambda x(\lambda x(zx))$  is  $\alpha$ -equivalent to  $\lambda x(\lambda y(zy))$ , where these later terms obey the Barendregt variable convention. The proof of the following theorem is also left to the exercises.

**Theorem 1.39**

For every  $\lambda$ -term  $E$ , there is a  $\lambda$ -term  $E'$  which is  $\alpha$ -equivalent to  $E$  such that  $E'$  obeys the Barendregt variable convention.

If then instead of typing  $\lambda x(\lambda x(x))$  we type the  $\alpha$ -equivalent term  $\lambda x(\lambda y(y))$  that obeys the Barendregt variable convention, we have the derivation:

$$\frac{\frac{x : \sigma, y : \tau \vdash y : \tau \quad (\text{Var})}{x : \sigma \vdash \lambda y(y) : \tau \rightarrow \tau} (\text{Abs})}{\vdash \lambda x(\lambda y(y)) : \sigma \rightarrow (\tau \rightarrow \tau)} (\text{Abs})$$

where no silent replacement of terms with  $\alpha$ -equivalent terms is necessary. We will in fact shortly see a theorem one of whose consequences is that silent replacement of terms with  $\alpha$ -equivalent terms can be avoided in typing derivations for  $\lambda$ -terms that obey the Barendregt variable convention.

In all our examples so far, we have proved that  $\tau$  is a possible type of  $E$  by constructing a proof of the sequent  $\vdash E : \tau$ . That is to say, we have proved that in the *empty* context  $E$  can be assigned type  $\tau$ . Is it always true that when  $E$  can be assigned type  $\tau$ , we can prove  $\vdash E : \tau$  in our system?

The answer is clearly no. Consider an *open*  $\lambda$ -term such as the variable  $x$ . This term can be given any type  $\tau$  we please. However, for no  $\tau$  can we prove  $\vdash x : \tau$ . For it is clear that the final step in such a proof could not be (App) or (Abs), as these rules produce conclusions  $\Gamma \vdash E : \tau$  in which  $E$  is *complex* – in any application of (App),  $E$  will have the form  $E_1(E_2)$ , and in any application of (Abs),  $E$  will have the form  $\lambda v(E)$  – while the variable  $x$  is *simple*. The last step of the proof must therefore be a leaf of the form (Var). However all leaves in  $\mathbf{TR}_0$  have non-empty contexts and thus a leaf on its own cannot constitute a proof of  $\vdash x : \tau$ .

On the other hand, the proof consisting of nothing other than the leaf  $x : \tau \vdash x : \tau$  shows that in some context (namely, the context  $x : \tau$ ),  $x$  has the possible type  $\tau$ . The following theorem tells us more generally that typing derivations for untyped *open*  $\lambda$ -terms always involve non-empty contexts in which the type of each free variable is declared

**Theorem 1.40**

In any proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0$ , for every free variable  $v$  of  $E$  the context  $\Gamma$  contains a declaration of the form  $v : \tau$  for some  $\tau$ .

The proof is a straightforward induction on the construction of typing derivation in  $\mathbf{TR}_0$  and is left to the exercises. As a kind of converse of this result, it can also easily be shown that if there is any typing derivation whose conclusion is of the form  $\Gamma \vdash E : X$ , then there is a typing derivation whose conclusion has the form  $\Gamma' \vdash E : X$ , where the context  $\Gamma'$  consists of nothing but variable declarations for the free variables of  $E$ . In fact, we will show something stronger. First, consider



the following theorem, which tells us that typing derivations for terms that obey the Barendregt variable convention may be assumed to have a very well behaved structure.

#### Theorem 1.41

Suppose  $E$  is a  $\lambda$ -term obeying the Barendregt variable convention, and that  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0$ . Then there is a proof of  $\Gamma' \vdash E : X$  in  $\mathbf{TR}_0$  for some  $\Gamma' \subseteq \Gamma$  such that (i) there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, (ii) the variables whose types are declared in  $\Gamma'$  are precisely the free variables of  $E$ , (iii) all the elements of  $\Gamma'$  are distinct, and (iv) every variable declaration in the proof is either an element of  $\Gamma'$ , or has the form  $v : \tau$  where  $v$  is a bound variable of  $E$ .

The proof of this theorem is again a straightforward induction on the typing derivations of  $\mathbf{TR}_0$  and is left to the exercises. We then have the following corollary.

#### Corollary 1.42

Suppose there is a typing derivation of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0$  (where  $E$  does not necessarily obey the Barendregt variable convention.) Then there is a typing derivation  $\Gamma' \vdash E : X$  in  $\mathbf{TR}_0$  such that  $\Gamma'$  has the form  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ , where  $v_1, \dots, v_n$  are the free variables of  $E$ . (Note that there is no repetition.) Moreover, in this typing derivation there is at most one silent replacement of a  $\lambda$ -term by an  $\alpha$ -equivalent term, and it occurs in the last step.

#### Proof

Suppose we have a typing derivation  $T_0$  whose conclusion is a sequent  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0$ , and let  $E'$  be a term  $\alpha$ -equivalent to  $E$  that obeys the Barendregt variable convention. By silently replacing  $E$  with the  $\alpha$ -equivalent term  $E'$  in the last step of  $T_0$ , we can construct a typing derivation  $T_1$  of  $\Gamma \vdash E' : X$ . Applying Theorem 1.41, we then have a typing derivation  $T_2$  of  $\Gamma' \vdash E' : X$  such that  $\Gamma'$  is a list of typing declarations for all the free variables of  $E'$  (and no others) without repetition, and such that there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent  $\lambda$ -terms in  $T_2$ . Note that the free variables of  $E$  and  $E'$  are the same. If in the last step of  $T_2$  we introduce a silent replacement of  $E'$  with  $E$ , we then have a typing derivation  $T_3$  whose conclusion is  $\Gamma' \vdash E : X$  in  $\mathbf{TR}_0$ . In  $T_3$ ,  $\Gamma'$  consists only of typing declarations for all the free variables of  $E$  (and no others) without repetition. There is also at most one silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, and it occurs in the last step.

This tells us that even typing derivations for terms that do not obey the Barendregt variable convention may be assumed to have a well behaved structure.

We now turn to the question of in what sense  $\mathbf{TR}_0$  allows us to prove all the basic facts about typing that we might want to prove. Indeed, there is a sense in which  $\mathbf{TR}_0$  is sound and complete. This is captured by the following theorem:

**Theorem 1.43: Soundness and Completeness of  $\mathbf{TR}_0$ .**

Suppose  $E$  is an untyped  $\lambda_0$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be simple types. Then the following are equivalent:

- (i) It is possible to assign simple types (i.e., types of  $\lambda_0^{\text{type}}$ ) to all the variables of  $E$  in such a way that the free variables  $v_1, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \dots, \tau_n$ , and  $E$  itself has type  $X$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

A consequence of this theorem is that  $\mathbf{TR}_0$  proves all true claims of the form

$$E \text{ has } X \text{ as a possible type in context } \Gamma$$

for any untyped  $\lambda$ -term  $E$ , type  $X$ , and context  $\Gamma$  in which types are assigned to all the free variables of  $E$ . In this sense, the system  $\mathbf{TR}_0$  proves every basic fact we might want to know about typing.

Unlike many soundness and completeness proofs in logic, the proof of Theorem 1.43 is fairly straightforward. That (i) entails (ii) (the ‘completeness’ of  $\mathbf{TR}_0$ ) can be proven by an easy induction on the construction of untyped  $\lambda$ -terms  $E$ , and uses Theorem 1.41. That (ii) entails (i) (the ‘soundness’ of  $\mathbf{TR}_0$ ) can be proven by an easy induction on the construction of proofs of  $\mathbf{TR}_0$ . The proof is left to the appendix, but you should try proving both directions on your own before consulting it.

A final important result is the following:

**Theorem 1.44: Subject Reduction Theorem for  $\mathbf{TR}_0$ .**

In  $\mathbf{TR}_0$ , if  $\Gamma \vdash M : \sigma$  and  $M \rightarrow_{\beta\eta} N$ , then  $\Gamma \vdash N : \sigma$ .

This theorem tells us that neither  $\beta$ -reduction nor  $\eta$ -reduction decreases the set of possible types an untyped term may be assigned in the Curry-style typed  $\lambda$ -calculus. The proof of this theorem is also left to the exercises. Note, however, that reduction can *increase* the set of possible types an untyped term can be assigned. You will also demonstrate this in an exercise.

The Curry-style typed  $\lambda$ -calculus is not really a fundamentally different sort of thing from the Church-style typed  $\lambda$ -calculus. In both cases, we are ultimately concerned with the way in which untyped  $\lambda$ -terms (that is, terms of  $\lambda_0$ ) may be typed (that is, transformed into terms of  $\lambda_0^{\text{type}}$ .) Church and Curry simply used different formalisms to analyse this basic phenomenon. Both formalisms will be useful to us in what follows.

**Exercises for Section 1.12**

1. Construct typing derivations in  $\mathbf{TR}_0$  that prove the following.

- (i)  $\vdash \lambda x(x(\lambda y(y))) : ((\sigma \rightarrow \sigma) \rightarrow \tau) \rightarrow \tau$
- (ii)  $\vdash \lambda x \lambda y(y(x(y))) : ((\tau \rightarrow \sigma) \rightarrow \tau) \rightarrow ((\tau \rightarrow \sigma) \rightarrow \sigma)$
- (iii)  $\vdash \lambda x \lambda y \lambda z(x(z)(y(z))) : (\sigma \rightarrow (\tau \rightarrow \rho)) \rightarrow ((\sigma \rightarrow \tau) \rightarrow (\sigma \rightarrow \rho))$ .

2. Construct a typing derivation in  $\mathbf{TR}_0$  proving a conclusion of the form

$$\vdash \lambda x(x)(\lambda y(y)(\lambda z(z))) : \tau$$

for some  $\tau$ . Then construct a typing derivation of  $\vdash s : \tau$  for the same  $\tau$ , where  $s$  is the  $\beta$ -normal form of  $\lambda x(x)(\lambda y(y)(\lambda z(z)))$ .

3. Construct an untyped closed  $\lambda$ -term  $\alpha$  such that  $\alpha$  cannot be assigned any type in the Curry-style  $\lambda$ -calculus (and thus there is no proof of  $\vdash \alpha : \tau$  for any  $\tau$ ), but for some untyped  $\lambda$ -term  $\alpha'$  with  $\alpha \rightarrow_{\beta} \alpha'$  the term  $\alpha'$  can be assigned a type in the Curry-style  $\lambda$ -calculus (that is, there is some proof of  $\vdash \alpha' : \sigma$  for some  $\sigma$ ).

4. Construct an untyped closed  $\lambda$ -term  $\alpha$  such that  $\alpha$  can be assigned a type in the Curry-style  $\lambda$ -calculus, but for some untyped  $\lambda$ -term  $\alpha'$  with  $\alpha \rightarrow_{\beta} \alpha'$  the term  $\alpha'$  can be assigned a type in the Curry-style  $\lambda$ -calculus that  $\alpha$  cannot be assigned.

5. Prove Theorem 1.38.

6. Prove Theorem 1.39.

7. Prove Theorem 1.40.

8. Prove Theorem 1.41.

9. Are the following true or false? If true, provide a proof. If false, provide a counterexample and explain why it is a counterexample.

- (i) If  $\Gamma \vdash E : X$ ,  $v$  is a variable not declared in  $\Gamma$ , and  $Y$  is a type, then  $\Gamma, v : Y \vdash E : X$ . (Sometimes this is called *weakening*.)
- (ii) If  $\Gamma, v : Y, v : Y \vdash E : X$ , then  $\Gamma, v : Y \vdash E : X$ . (Sometimes this is called *contraction*.)
- (iii) If  $\Gamma \vdash E : X$ , then every variable declaration in the proof is either an element of  $\Gamma$ , or has the form  $v : \tau$  where  $v$  is a bound variable of  $E$ .
- (iv) If  $\Gamma \vdash E : X$  and there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, then  $E$  obeys the Barendregt variable convention.
- (v) If  $\Gamma \vdash E : X$  and there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, then every variable declaration in the proof is either an element of  $\Gamma$ , or has the form  $v : \tau$  where  $v$  is a bound variable of  $E$ .

10. Prove Theorem 1.44. (Hint: use the Soundness and Completeness Theorem to reduce the

problem to that of showing that if  $M \rightarrow_{\beta} N$ , and if consistent with the free variables of  $M$  being assigned certain types,  $M$  itself may be assigned the type  $\sigma$ , then consistent with the free variables of  $N$  being assigned those same types,  $N$  itself may also be assigned the type  $\sigma$ .)

### 1.13 The Bound-typed $\lambda$ -calculus.

In this section, we describe a second version of the Curry-style  $\lambda$ -calculus that will be used very frequently in the remainder of the book.

One slight oddity of the terminology

$$\Gamma \vdash E : X$$

is that it only expresses the idea that under the typing assumptions in  $\Gamma$ , the  $\lambda$ -term  $E$  *may* be given the type  $X$ . This does not mean that  $E$  *must* be given the type  $X$ . For example, for any type  $A$  we can easily derive in  $\mathbf{TR}_0$  that

$$\vdash \lambda x(x) : A \rightarrow A.$$

For each type  $A$ , this means that the type  $A \rightarrow A$  is a type that  $\lambda x(x)$  *may* receive. For no particular  $A$  is  $A \rightarrow A$  the type that  $\lambda x(x)$  *must* receive.

However, typically in logic the relation  $\Gamma \vdash S$  denotes a type of *necessitation*, and says that if  $\Gamma$  is true, then  $S$  *must* be true. We can bring our notation  $\Gamma \vdash E : X$  into conformity with this idea of necessitation by moving to a so-called *partially typed*  $\lambda$ -calculus, which we now explain.

Consider the following definition

**Definition 1.45: Bound-typed  $\lambda$ -terms.**

A  $\lambda$ -term is *bound-typed* if all its bound variables are typed, and all its free variables are untyped.

So for example, the following terms are all bound-typed:

$$\begin{aligned} &\lambda x^{\alpha}(x^{\alpha}), \\ &\lambda x^{\alpha}(yx^{\alpha}), \\ &\lambda x^{\alpha}(x^{\alpha}(\lambda z^{\beta}(yz^{\beta}))). \end{aligned}$$

Using our former conventions, when there is no ambiguity we can eliminate the type symbols on bound variables except when they are ‘declared’ immediately next to a  $\lambda$  symbol. For example, we can write the above three expressions more compactly as follows:

$$\begin{aligned} &\lambda x^{\alpha}(x), \\ &\lambda x^{\alpha}(yx), \end{aligned}$$

$$\lambda x^\alpha(x(\lambda z^\beta(yz))).$$

Some bound-typed  $\lambda$ -terms can have further superscripts added to them and be transformed to a full  $\lambda$ -term of  $\lambda_0^{\text{type}}$ . In this way, such terms then acquire an overall type. For example, consider the term  $\lambda x^\alpha(yx)$  (that is,  $\lambda x^\alpha(yx^\alpha)$ .) Given the occurrence of the subexpression  $yx^\alpha$ , for this to be transformed into a grammatical  $\lambda$ -term of  $\lambda_0^{\text{type}}$  the variable  $y$  must have a type of the form  $\alpha \rightarrow \beta$ . Assuming then that  $y$  has type  $\alpha \rightarrow \beta$ ,  $yx^\alpha$  has type  $\beta$ , and so the bound-typed  $\lambda$ -term  $\lambda x^\alpha(yx^\alpha)$  is transformed into the  $\lambda_0^{\text{type}}$  term  $\lambda x^\alpha(y^{\alpha \rightarrow \beta}x^\alpha)$  with overall type  $\alpha \rightarrow \beta$ . The set of possible types of  $\lambda x^\alpha(yx^\alpha)$  is therefore the set of types of the form  $\alpha \rightarrow \beta$ .

As before, we will say that bound-typed  $\lambda$ -terms that are identical up to a renaming of bound variables are  $\alpha$ -equivalent. (All the delicate provisos discussed earlier still apply.) When renaming bound variables, any types associated with such variables remain fixed. So for example, the terms

$$\lambda x^\alpha(x) \quad \lambda y^\alpha(y)$$

are  $\alpha$ -equivalent, but

$$\lambda x^\alpha(x) \quad \lambda y^\beta(y)$$

are not.

In the version of the Curry-style typed  $\lambda$ -calculus we develop in this section, we will be assigning types to bound-typed terms. We do so by deriving sequents of the form

$$v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n \vdash E : \sigma \tag{19}$$

where  $E$  is a bound-typed term and the list of variables  $v_1, v_2, \dots, v_n$  includes (but is not necessarily limited to) all the free variables of  $E$ . As before, this sequent will have the interpretation

Consistent with the constraint that the variables  $v_1, v_2, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, it is possible to type  $E$  so that it becomes a  $\lambda_0^{\text{type}}$  term of type  $\sigma$ .

Note however that the type of a term is completely determined by specifying the types of all its variables (both free and bound). Because in a bound-typed term the types of bound variables are already given, it follows that the type of a bound-typed term is uniquely determined by the types of its free variables. This means, however, that

$$v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n \vdash E : \sigma$$

will be true for at most one  $\sigma$ . Thus we may interpret (19) as saying

If the variables  $v_1, v_2, \dots, v_n$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, then  $E$  *must* be assigned type  $\sigma$ .

In this way, our type-theoretic sequents now make a claim about *necessitation*, that is, what *must* be the case if something else is the case.

In the version of the  $\lambda$ -calculus we consider here, we will be able to derive sequents such as

$$y : \alpha \rightarrow \beta \vdash \lambda x^\alpha(yx) : \alpha \rightarrow \beta.$$

This sequent may be understood as saying that if in  $\lambda x^\alpha(yx)$  the free variable  $y$  gets assigned the type  $\alpha \rightarrow \beta$ , then the entire term  $\lambda x^\alpha(yx)$  *must* get assigned the type  $\alpha \rightarrow \beta$ .

We will also allow sequents such as

$$y : \alpha \rightarrow \beta, u : \gamma \vdash \lambda x^\alpha(yx) : \alpha \rightarrow \beta.$$

This says that if in  $\lambda x^\alpha(yx)$  any free occurrence of  $y$  gets assigned the type  $\alpha \rightarrow \beta$  and any free occurrence of  $u$  gets assigned the type  $\gamma$ , then the entire term  $\lambda x^\alpha(yx)$  *must* get assigned the type  $\alpha \rightarrow \beta$ . This is trivially true because, in addition to the previous considerations,  $u$  does not actually occur free in the  $\lambda$ -term in question.

We furthermore even allow sequents such as

$$y : \alpha \rightarrow \beta, x : \gamma \vdash \lambda x^\alpha(yx) : \alpha \rightarrow \beta.$$

Such a sequent looks confusing at first, as  $x$  appears to be typed differently on the left and right. However as before we will allow ‘silent’ replacing of  $\lambda$ -terms with  $\alpha$ -equivalent  $\lambda$ -terms. This means that the sequent just written can be re-expressed in a less confusing way as follows

$$y : \alpha \rightarrow \beta, x : \gamma \vdash \lambda u^\alpha(yu) : \alpha \rightarrow \beta,$$

where we have replaced the  $\lambda$ -term  $\lambda x^\alpha(yx)$  with the  $\alpha$ -equivalent term  $\lambda u^\alpha(yu)$ . This later sequent is of course entirely unproblematic.

Let us think about the rules that our new system of sequents should obey. First of all, as before we can keep the (Var) rule

$$\Gamma, v : \tau \vdash v : \tau \quad (\text{Var})$$

where  $v$  is a variable. The variable  $v$  of course a bound-typed  $\lambda$ -term, and thus this axiom just makes the trivial claim that if among other variable declarations  $v$  has been given type  $\tau$ , then the bound-typed term  $v$  must be given type  $\tau$ .

Next, the (App) rule can also be used in its familiar form

$$\frac{\Gamma \vdash E_1 : \tau \rightarrow \sigma \quad \Gamma \vdash E_2 : \tau}{\Gamma \vdash E_1(E_2) : \sigma} \quad (\text{App})$$

where  $E_1$  and  $E_2$  are bound-typed terms. For if in context  $\Gamma$  the bound-typed terms  $E_1$  and  $E_2$  must be assigned types  $\tau \rightarrow \sigma$  and  $\tau$  respectively, then in this same context  $\Gamma$  the term  $E_1(E_2)$  (which will also be bound-typed) must be assigned type  $\sigma$ .

The (Abs) rule will clearly require a small modification, as follows

$$\frac{\Gamma, v : \tau \vdash E : \sigma}{\Gamma \vdash \lambda v^\tau(E) : \tau \rightarrow \sigma} \quad (\text{Abs}^{\text{bt}})$$

Here  $E$  is a bound-typed term. The only difference between this and the previous version of the (Abs) rule is that in the conclusion, the bound variable  $v$  in the  $\lambda$ -term is explicitly assigned type  $\tau$ . (We indicate this difference with the superscript **bt** in  $\text{Abs}^{\text{bt}}$ , indicating that it is the ‘bound-typed’ variant of the rule.) Suppose that as stated in the premise of this rule, in the context  $\Gamma$  and under the assumption that the variable  $v$  has type  $\tau$ , the bound-typed term  $E$  must receive the type  $\sigma$ . This then means that in the context  $\Gamma$ , the bound-typed term  $\lambda v^\tau(E)$  must receive the type  $\tau \rightarrow \sigma$ .

In sum, we have the following typing rules for our system:

**Definition 1.46: Rules of the bound-typed  $\lambda$ -calculus  $\text{TR}_0^{\text{bt}}$ .**

- (i)  $\Gamma, v : \tau \vdash v : \tau$  (Var)
- (ii) 
$$\frac{\Gamma \vdash E_1 : \tau \rightarrow \sigma \quad \Gamma \vdash E_2 : \tau}{\Gamma \vdash E_1(E_2) : \sigma}$$
 (App)
- (iii) 
$$\frac{\Gamma, v : \tau \vdash E : \sigma}{\Gamma \vdash \lambda v^\tau(E) : \tau \rightarrow \sigma}$$
 ( $\text{Abs}^{\text{bt}}$ )

where  $v$  is a variable,  $E_1$ ,  $E_2$  and  $E$  are bound-typed terms of the  $\lambda$ -calculus,  $\sigma, \tau$  are types, and all contexts are assumed to be grammatical.

We call this set of rules  $\text{TR}_0^{\text{bt}}$ . It is of course simply a trivial variant of  $\text{TR}_0$ , with the obvious types attached to bound variables in  $\lambda$ -terms where needed.

Even though their interpretations are different, any proof in  $\text{TR}_0$  can be modified easily to give a proof in  $\text{TR}_0^{\text{bt}}$ . For example, the following proof in  $\text{TR}_0$ :

$$\frac{\frac{\frac{x : \sigma, y : \sigma \rightarrow \tau \vdash x : \sigma \quad x : \sigma, y : \sigma \rightarrow \tau \vdash y : \sigma \rightarrow \tau}{x : \sigma, y : \sigma \rightarrow \tau \vdash yx : \tau} \text{ (App)}}{x : \sigma \vdash \lambda y(yx) : (\sigma \rightarrow \tau) \rightarrow \tau} \text{ (Abs)}}{\vdash \lambda x(\lambda y(yx)) : \sigma \rightarrow ((\sigma \rightarrow \tau) \rightarrow \tau)} \text{ (Abs)}$$

can be transformed into the following proof in  $\text{TR}_0^{\text{bt}}$ :

$$\frac{\frac{\frac{x : \sigma, y : \sigma \rightarrow \tau \vdash x : \sigma \quad x : \sigma, y : \sigma \rightarrow \tau \vdash y : \sigma \rightarrow \tau}{x : \sigma, y : \sigma \rightarrow \tau \vdash yx : \tau} \text{ (App)}}{x : \sigma \vdash \lambda y^{\sigma \rightarrow \tau}(yx) : (\sigma \rightarrow \tau) \rightarrow \tau} \text{ (Abs}^{\text{bt}}\text{)}}{\vdash \lambda x^\sigma(\lambda y^{\sigma \rightarrow \tau}(yx)) : \sigma \rightarrow ((\sigma \rightarrow \tau) \rightarrow \tau)} \text{ (Abs}^{\text{bt}}\text{)}$$

by typing the bound variables in the  $\lambda$ -terms in the obvious way. One can also go from the  $\text{TR}_0^{\text{bt}}$  proof to the original  $\text{TR}_0$  proof by simply removing all types from bound variables. So the structure of proofs in  $\text{TR}_0$  and  $\text{TR}_0^{\text{bt}}$  are for all intents and purposes the same.

As a result, for many theorems of  $\text{TR}_0$  there is a corresponding theorem about  $\text{TR}_0^{\text{bt}}$ . For example, we have the following.

**Theorem 1.47**

In any proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0^{\text{bt}}$ , for every free variable  $v$  of  $E$  the context  $\Gamma$  contains a declaration of the form  $v : \tau$  for some  $\tau$ .

**Theorem 1.48**

Suppose  $E$  is a bound-typed  $\lambda$ -term obeying the Barendregt variable convention, and that  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0^{\text{bt}}$ . Then there is a proof of  $\Gamma' \vdash E : X$  in  $\mathbf{TR}_0^{\text{bt}}$  for some  $\Gamma' \subseteq \Gamma$  such that (i) there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, (ii) the variables whose types are declared in  $\Gamma'$  are precisely the free variables of  $E$ , (iii) all the elements of  $\Gamma'$  are distinct, and (iv) every variable declaration in the proof is either an element of  $\Gamma'$ , or has the form  $v : \tau$  where  $v$  is a bound variable of  $E$ .

**Theorem 1.49: Soundness and Completeness of  $\mathbf{TR}_0^{\text{bt}}$ .**

Suppose  $E$  is a bound-typed  $\lambda$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) If the variables  $v_1, v_2, \dots, v_n$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, then  $E$  must be assigned type  $\sigma$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_0^{\text{bt}}$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

**Theorem 1.50: Subject Reduction Theorem for  $\mathbf{TR}_0^{\text{bt}}$ .**

In  $\mathbf{TR}_0^{\text{bt}}$ , if  $\Gamma \vdash M : \sigma$  and  $M \rightarrow_{\beta\eta} N$ , then  $\Gamma \vdash N : \sigma$ .

These proofs may all be obtained by minor modifications of the proofs of the corresponding theorems for  $\mathbf{TR}_0$ . Verifying this is left to the exercises and the appendix to this chapter.

In the chapters that follow, we will be interested in expanding the system  $\mathbf{TR}_0$  to richer systems  $\mathbf{TR}_1, \mathbf{TR}_2, \dots$  in which we reason in more complicated ways about the typing of untyped  $\lambda$ -terms. As in the case of  $\mathbf{TR}_0$ , trivial modifications of these systems will yield richer systems  $\mathbf{TR}_1^{\text{bt}}, \mathbf{TR}_2^{\text{bt}}, \dots$  in which we reason in more complicated ways about the typing of bound-typed  $\lambda$ -terms.



**Exercises for Section 1.13**

1. Construct typing derivations in  $\mathbf{TR}_0^{\text{bt}}$  that prove the following.

(i)  $\vdash \lambda x^\sigma (\lambda y^\tau (x)) : \sigma \rightarrow (\tau \rightarrow \sigma)$

(ii)  $\vdash \lambda x^{\sigma \rightarrow (\tau \rightarrow \rho)} \lambda y^{\sigma \rightarrow \tau} \lambda z^\sigma (x(z)(y(z))) : (\sigma \rightarrow (\tau \rightarrow \rho)) \rightarrow ((\sigma \rightarrow \tau) \rightarrow (\sigma \rightarrow \rho))$ .

2. Prove Theorem 1.47.

3. Prove Theorem 1.48.

4. Prove Theorem 1.50.

**1.14 Other Resources**

For a basic and very readable introduction to the  $\lambda$ -calculus, see chapters 1 through 4 of Selinger's *Lecture Notes on the Lambda Calculus* [9]. For a more extended treatment, see Hindley and Seldin's *Lambda-Calculus and Combinators* [4]. For an extremely detailed treatment, see Barendregt's *The Lambda Calculus, Its Syntax and Semantics* [1]. Chapters 1 and 2 of Hindley's *Basic Simple Type Theory* [5] also contain a more condensed but very good introduction to the untyped  $\lambda$ -calculus, and chapters 1 and 3 of Sorenson and Urzyczyn's *Lectures on the Curry-Howard Isomorphism* [10] contains a very good introduction to both the untyped and typed  $\lambda$ -calculus.

For a more general discussion of the themes of the chapter, Barendregt's *The Lambda Calculus, Its Syntax and Semantics* [2] is also worth reading.

**1.15 Appendix**

In this section, we provide proofs for several results in the main text.

**Theorem 1.13: The Church-Rosser Theorem for  $\beta$ -reductions in  $\lambda_0$ .**

In the untyped  $\lambda$ -calculus  $\lambda_0$ , if  $X \rightarrow_\beta Y_1$  and  $X \rightarrow_\beta Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_\beta Z$  and  $Y_2 \rightarrow_\beta Z$ .

**Proof of Theorem 1.13.**

In this proof, all  $\lambda$ -terms are understood to be terms of the untyped  $\lambda$ -calculus  $\lambda_0$ . Without loss of generality, we assume all terms obey the Barendregt variable convention.

It is useful to begin by defining a more general notion of reduction. Given a  $\lambda$ -term  $e$ , there may be many distinct 1-step  $\beta$ -reductions that could potentially be performed on it,

each corresponding to a different redex of the form  $\lambda v(t_1)(t_2)$ . Given any two distinct redexes  $r = \lambda v(t_1)(t_2)$  and  $r' = \lambda v'(t'_1)(t'_2)$ , they will either be disjoint,  $r$  will be a subexpression of  $r'$ , or  $r'$  will be a subexpression of  $r$ . The set of redexes contained in  $e$  thus form a finite set of trees under the relation of inclusion which we call the *redex trees* of  $e$ ; we assume these trees are oriented so that redexes containing no other redexes are at the top.

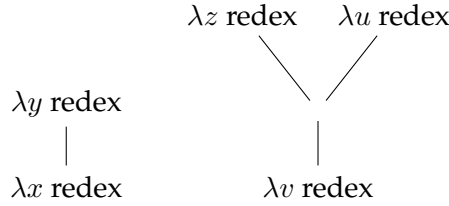
For example, if  $e$  is

$$\lambda x(x)(\lambda y(y)(y)) w(\lambda v(v(\lambda z(z)(z)))(\lambda w(\lambda u(u)(u)))) \quad (20)$$

$\underbrace{\lambda x(x)(\lambda y(y)(y))}_{\lambda x \text{ redex}} \quad \underbrace{w(\lambda v(v(\lambda z(z)(z)))(\lambda w(\lambda u(u)(u))))}_{\lambda v \text{ redex}}$

$\lambda y \text{ redex} \quad \lambda z \text{ redex} \quad \lambda u \text{ redex}$

then there are 5 redexes in  $e$  as depicted. The redex trees of  $e$  are then as follows



Note that if one redex  $r'$  lies above another redex  $r = \lambda v(t_1)(t_2)$  in the redex trees, then not only is  $r'$  a subexpression of  $r$ , but it is in fact a (not necessarily proper) subexpression of  $t_1$  or  $t_2$ . This is because  $r'$  is a proper subexpression of  $r = \lambda v(t_1)(t_2)$ , and the only proper subexpressions of  $\lambda v(t_1)(t_2)$  are  $\lambda v(t_1)$  and the (not necessarily proper) subexpressions of  $t_1$  and  $t_2$ . The term  $\lambda v(t_1)$  does not have the form of a redex, and thus  $r'$  must be a (not necessarily proper) subexpression of  $t_1$  or  $t_2$ .

As a consequence of this, if  $r'$  lies above  $r$  and a 1-step  $\beta$ -reduction is performed on  $r'$ , the redex  $r = \lambda v(t_1)(t_2)$  is transformed into a term  $\lambda v(t_1^*)(t_2)$  or  $\lambda v(t_1)(t_2^*)$ . Either way, this still has the form of a redex. Because of this, it is possible to perform  $\beta$ -reduction on as many (or as few) of the redexes contained in the redex trees of  $e$  that we want, so long as we work from the innermost redexes outwards – that is, so long as we reduce redexes starting with those uppermost in the redex trees, then continue reducing other redexes in the redex trees of  $e$ , never moving upwards. If doing so produces a term  $e'$ , we will say that  $e \Rightarrow_{\beta} e'$ . If we reduce *all* the redexes contained in  $e$  (working from the uppermost leaves of the redex trees and moving downwards), we call the resulting term  $e^*$  (or when there is risk of confusion,  $[e]^*$ ), which we call the *complete reduction* of  $e$ . (Note that the order in which we perform  $\beta$ -reductions associated with disjoint redexes does not matter.) By definition, we then have  $e \Rightarrow_{\beta} e^*$ . For example, for the term  $e^*$  given in (20), the complete reduction of  $e$  is given by  $e^* = yw(\lambda w(u)(z))$ . (Note that this term is not necessarily in  $\beta$ -normal form.)

We then have the following sequence of lemmas.

**Lemma 1:** For any  $\lambda$ -terms  $e$  and  $e'$ ,  $e \rightarrow_\beta e'$  entails  $e \Rightarrow_\beta e'$ .

**Proof:** The reduction  $e \rightarrow_\beta e'$  corresponds to the case in which we perform a  $\beta$ -reduction on only one redex of  $e$ , and so  $e \Rightarrow_\beta e'$  is immediate from the definition of  $\Rightarrow_\beta$ .  $\square$

**Lemma 2:** For any  $\lambda$ -terms  $e$  and  $e'$ ,  $e \Rightarrow_\beta e'$  entails  $e \twoheadrightarrow_\beta e'$ .

**Proof:** If  $e \Rightarrow_\beta e'$ , the reduction of  $e$  to  $e'$  may be thought of as a sequence of  $\beta$ -reductions (never moving upwards in the redex trees of  $e$ ), and thus  $e \twoheadrightarrow_\beta e'$ .

**Lemma 3:** For any  $\lambda$ -terms  $s, s', t, t'$ , if  $s \Rightarrow_\beta s'$  and  $t \Rightarrow_\beta t'$ , then  $s[t/v] \Rightarrow_\beta s'[t'/v]$ .

**Proof:** Assume  $s \Rightarrow_\beta s'$  and  $t \Rightarrow_\beta t'$ . The term  $s[t/v]$  may be  $\beta$ -reduced to the term  $s'[t'/v]$  as follows. First reduce every occurrence of  $t$  that appears as a result for substituting  $t$  for  $v$  in  $s[t/v]$  to  $t'$  (using  $t \Rightarrow_\beta t'$ ), thereby obtaining  $s[t'/v]$ . Then perform the same reductions in  $s \Rightarrow_\beta s'$  to obtain a  $\beta$ -reduction from  $s[t'/v]$  to  $s'[t'/v]$ . (The fact that we have substituted  $t'$  for the free variable  $v$  has no effect on this sequence of reductions.) The result is a sequence of  $\beta$ -reductions from  $s[t/v]$  to  $s'[t'/v]$  which is easily seen to never move upwards in the redex trees of  $s[t/v]$ . Thus  $s[t/v] \Rightarrow_\beta s'[t'/v]$ .  $\square$

**Lemma 4:** For any  $\lambda$ -terms  $s$  and  $t$ , if  $s \Rightarrow_\beta t$  then  $t \Rightarrow_\beta s^*$ .

For a  $\lambda$ -term  $s$ , let  $\text{height}(s)$  be the length of the longest branch in the redex trees of  $s$ . (For example, if  $e$  is the term given in (20), then  $\text{height}(e) = 2$ .) The proof of the lemma is then by induction on  $\text{height}(s)$ .

If  $\text{height}(s) = 0$ , this means that  $s$  contains no redexes, and so  $s^* = s$ . From  $s \Rightarrow_\beta t$ , we must also have  $t = s$ . Thus  $t \Rightarrow_\beta s^*$  holds iff  $s \Rightarrow_\beta s$  holds, which is immediate.

Suppose now that the result holds for terms  $s'$  with  $\text{height}(s') \leq l$ , and consider a term  $s$  with  $\text{height}(s) = l + 1$ . We begin by considering the case in which  $s$  itself is a redex, and so has the form

$$s = \lambda x(M)(N),$$

with  $\text{height}(M), \text{height}(N) \leq l$ . Note that we have

$$s^* = M^*[N^*/x],$$

because the complete reduction of  $\lambda x(M)(N)$  is obtained by completely  $\beta$ -reducing  $M$  and  $N$ , and then performing the outmost and final  $\beta$ -reduction on the  $\lambda x$  redex.

Suppose next that  $s \Rightarrow_\beta t$ . There are two cases to consider

**(Case 1)** In the reduction  $s \Rightarrow_\beta t$ , the redex  $s$  itself is among the redexes that are reduced.

**(Case 2)** In the reduction  $s \Rightarrow_\beta t$ , the redex  $s$  itself is not among the redexes that are reduced.

We consider these cases in turn. In Case 1, because any redex properly contained in  $\lambda x(M)(N)$  must either be contained in  $M$  or contained in  $N$ , we know that  $t$  has the form

$$t = M'[N'/x] \quad \text{where } M \Rightarrow_\beta M' \text{ and } N \Rightarrow_\beta N'.$$

By inductive hypothesis, we have  $M' \Rightarrow_\beta M^*$  and  $N' \Rightarrow_\beta N^*$ . From  $t = M'[N'/x]$ ,  $M' \Rightarrow_\beta$

$M^*, N' \Rightarrow_\beta N^*$  and Lemma 3 we then have

$$t = M'[N'/x] \Rightarrow_\beta M^*[N^*/x] = s^*$$

so  $t \Rightarrow_\beta s^*$  as desired. Consider now Case 2. In this case,  $t$  has the form

$$t = \lambda x(M')(N') \quad \text{where } M \Rightarrow_\beta M' \text{ and } N \Rightarrow_\beta N'.$$

Again, by inductive hypothesis, we have  $M' \Rightarrow_\beta M^*$  and  $N' \Rightarrow_\beta N^*$ . From this we have

$$t = \lambda x(M')(N') \Rightarrow_\beta M^*[N^*/x] = s^*$$

and so  $t \Rightarrow_\beta s^*$  as desired. So either way,  $t \Rightarrow_\beta s^*$ . This completes the inductive step for the case in which  $s$  is a redex.

If  $s$  is not itself a redex and  $\text{height}(s) = l + 1$ , then  $s$  can be written

$$s = \dots r_1 \dots r_2 \dots r_i \dots$$

where the  $r_i$  are the redexes at the bottom of the redex trees of  $s$ , and for all  $i$ ,  $\text{height}(r_i) \leq l$ . Then we have

$$s^* = \dots r_1^* \dots r_2^* \dots r_i^* \dots$$

where the expressions occurring in  $\dots$  remain unchanged. If  $s \Rightarrow_\beta t$ , then  $t$  may be written

$$t = \dots r'_1 \dots r'_2 \dots r'_i \dots$$

where  $r_i \Rightarrow_\beta r'_i$  for each  $i$ . By inductive hypothesis we therefore have  $r'_i \Rightarrow_\beta r_i^*$ , and so

$$t = \dots r'_1 \dots r'_2 \dots r'_i \dots \Rightarrow_\beta \dots r_1^* \dots r_2^* \dots r_i^* \dots = s^*.$$

Thus  $t \Rightarrow_\beta s^*$  as desired. This completes the induction and the proof of the lemma.  $\square$

In what follows, let  $s^{[n]}$  be  $s^{**\dots*}$  with  $n$  asterisks, and let  $s^{[0]}$  be  $s$ . We then have

**Lemma 5:** If  $s_0 \Rightarrow_\beta s_1 \Rightarrow_\beta \dots \Rightarrow_\beta s_{n-1} \Rightarrow_\beta s_n$  then  $s_n^{[0]} \Rightarrow_\beta s_{n-1}^{[1]} \Rightarrow_\beta \dots \Rightarrow_\beta s_1^{[n-1]} \Rightarrow_\beta s_0^{[n]}$

**Proof:** By induction on  $n$ . If  $n = 1$ , then from  $s_0 \Rightarrow_\beta s_1$  and Lemma 4 we have  $s_1 \Rightarrow_\beta s_0^*$  as desired.

Suppose then that the lemma is true for  $n$ , and that

$$s_0 \Rightarrow_\beta s_1 \Rightarrow_\beta \dots \Rightarrow_\beta s_n \Rightarrow_\beta s_{n+1}. \quad (21)$$

By inductive hypothesis, from  $s_0 \Rightarrow_\beta s_1 \Rightarrow_\beta \dots \Rightarrow_\beta s_n$  we have

$$s_n^{[0]} \Rightarrow_\beta s_{n-1}^{[1]} \Rightarrow_\beta \dots \Rightarrow_\beta s_1^{[n-1]} \Rightarrow_\beta s_0^{[n]}. \quad (22)$$

We have  $s_n^{[0]} \Rightarrow_\beta s_{n-1}^{[1]}$  from (22) and  $s_n^{[0]} \Rightarrow_\beta s_{n+1}^{[0]}$  (i.e.,  $s_n \Rightarrow_\beta s_{n+1}$ ) from (21). Applying Lemma 4 to both of these, we have  $s_{n-1}^{[1]} \Rightarrow_\beta s_n^{[1]}$  and  $s_{n+1}^{[0]} \Rightarrow_\beta s_n^{[1]}$ . From  $s_{n-1}^{[1]} \Rightarrow_\beta s_{n-2}^{[2]}$  and  $s_{n-1}^{[1]} \Rightarrow_\beta s_n^{[1]}$  and Lemma 4 we have  $s_{n-2}^{[2]} \Rightarrow_\beta s_{n-1}^{[2]}$  and  $s_n^{[1]} \Rightarrow_\beta s_{n-1}^{[2]}$ . Continuing this reasoning, we end up with  $s_0^{[n]} \Rightarrow_\beta s_1^{[n]}$  and  $s_2^{[n-1]} \Rightarrow_\beta s_1^{[n]}$ . From  $s_0^{[n]} \Rightarrow_\beta s_1^{[n]}$  and Lemma 4 we then have  $s_1^{[n]} \Rightarrow_\beta s_0^{[n+1]}$ . This reasoning can be depicted more intuitively with the following diagram

$$\begin{array}{ccccccccc}
 s_n^{[0]} & \Longrightarrow & s_{n-1}^{[1]} & \Longrightarrow & s_{n-2}^{[2]} & \cdots & s_1^{[n-1]} & \Longrightarrow & s_0^{[n]} \\
 \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow \\
 s_{n+1}^{[0]} & \Longrightarrow & s_n^{[1]} & \Longrightarrow & s_{n-1}^{[2]} & \cdots & s_2^{[n-1]} & \Longrightarrow & s_1^{[n]} & \Longrightarrow & s_0^{[n+1]}
 \end{array}$$

Thus we have  $s_{n+1}^{[0]} \Rightarrow_\beta s_n^{[1]} \Rightarrow_\beta \dots \Rightarrow_\beta s_1^{[n]} \Rightarrow_\beta s_0^{[n+1]}$  as desired. This completes the proof of Lemma 5.  $\square$

To prove the Church-Rosser Theorem itself, suppose  $X \twoheadrightarrow_\beta Y_1$  and  $X \twoheadrightarrow_\beta Y_2$ . There are then sequences  $U_1, \dots, U_n$  and  $V_1, \dots, V_m$  of terms such that

$$X \twoheadrightarrow_\beta U_1 \twoheadrightarrow_\beta \dots \twoheadrightarrow_\beta U_n = Y_1 \quad \text{and} \quad X \twoheadrightarrow_\beta V_1 \twoheadrightarrow_\beta \dots \twoheadrightarrow_\beta V_m = Y_2$$

and so by Lemma 1 we have

$$X \Rightarrow_\beta U_1 \Rightarrow_\beta \dots \Rightarrow_\beta U_n = Y_1 \quad \text{and} \quad X \Rightarrow_\beta V_1 \Rightarrow_\beta \dots \Rightarrow_\beta V_m = Y_2.$$

Letting  $l$  be  $\max(n, m)$  and adding reductions of the form  $e \Rightarrow_\beta e$  to the shorter of these two sequences, we then have

$$X \Rightarrow_\beta U_1 \Rightarrow_\beta \dots \Rightarrow_\beta U_l = Y_1 \quad \text{and} \quad X \Rightarrow_\beta V_1 \Rightarrow_\beta \dots \Rightarrow_\beta V_l = Y_2.$$

From Lemma 5 it then follows that

$$Y_1 \Rightarrow_\beta \dots \Rightarrow_\beta X^{[l]} \quad \text{and} \quad Y_2 \Rightarrow_\beta \dots \Rightarrow_\beta X^{[l]}.$$

From Lemma 2 we then have  $Y_1 \twoheadrightarrow_\beta X^{[l]}$  and  $Y_2 \twoheadrightarrow_\beta X^{[l]}$ . Thus  $X^{[l]}$  is a term  $Z$  such that  $Y_1 \twoheadrightarrow_\beta Z$  and  $Y_2 \twoheadrightarrow_\beta Z$ .  $\square$

**Theorem 1.20: The Church-Rosser Theorem for  $\beta\eta$ -reductions in  $\lambda_0$ .**

In the untyped  $\lambda$ -calculus  $\lambda_0$ , if  $X \twoheadrightarrow_{\beta\eta} Y_1$  and  $X \twoheadrightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \twoheadrightarrow_{\beta\eta} Z$  and  $Y_2 \twoheadrightarrow_{\beta\eta} Z$ .

**Proof of Theorem 1.20.**

In this proof, all  $\lambda$ -terms are understood to be terms of the untyped  $\lambda$ -calculus  $\lambda_0$ . As before, we assume all terms obey the Barendregt variable convention. The proof follows the same strategy as the proof of Theorem 1.13, with which familiarity is assumed. We mostly restrict ourselves to just describing the way in which this earlier proof must be modified to produce a proof of Theorem 1.20.

Again, we define a more general notion of  $\beta\eta$ -reduction that we denote  $\Rightarrow_{\beta\eta}$ . As before, for a given  $\lambda$ -term  $e$  the set of  $\beta$ - or  $\eta$ -redexes contained in  $e$  form a finite set of trees when ordered under the relation of inclusion, with the innermost redexes lying at the top. We would like it to be the case that whenever one redex  $r$  lies higher than another redex  $\lambda x(M)(N)$  or  $\lambda x(E(x))$  (with  $x$  free in  $E$ ) that  $r$  is a subexpression of  $M$  or  $N$  or  $E$ , so that when a redex in a tree is reduced, the redexes below maintain their forms as redexes. But there are (precisely) two sorts of counterexamples to this:

$$\begin{array}{cc}
 \text{(a)} & \begin{array}{c} \lambda x(E(x)) \\ | \\ \lambda x(E(x))(N) \end{array} & \text{(b)} & \begin{array}{c} \lambda x(M)(z) \\ | \\ \lambda z(\lambda x(M)(z)) \end{array}
 \end{array}$$

- (a) the  $\eta$ -redex  $\lambda x(E(x))$  is a subexpression of of the  $\beta$ -redex  $\lambda x(E(x))(N)$  (where  $x$  is not free in  $E$ ), so that when  $\lambda x(E(x))$  is reduced to  $E$ ,  $\lambda x(E(x))(N)$  is reduced to  $E(N)$  (which need not be a redex)
- (b) the  $\beta$ -redex  $\lambda x(M)(z)$  is a subexpression of of the  $\eta$ -redex  $\lambda z(\lambda x(M)(z))$  (where  $z$  is not free in  $\lambda x(M)$ ), so that when  $\lambda x(M)(z)$  is reduced to  $M[z/x]$ ,  $\lambda z(\lambda x(M)(z))$  is reduced to  $\lambda z(M[z/x])$  (which need not be a redex)

To deal with this problem, we will declare any redex in  $e$  of the form  $\lambda x(E(x))(N)$  (where  $x$  is not free in  $E$ ) or  $\lambda z(\lambda x(M)(z))$  (where  $z$  is not free in  $\lambda x(M)$ ) to be *prohibited*. If we perform a 1-step reduction on any non-prohibited redex, then any other non-prohibited redex below it maintains its form as a redex. We can therefore perform a 1-step  $\beta$ - or  $\eta$ -reduction on as many (or as few) of the non-prohibited redexes contained in the redex trees of  $e$  that we want, working from the innermost redexes outwards – that is, working from the uppermost redexes downwards. As before, if doing so produces a term  $e'$ , we will say that  $e \Rightarrow_{\beta} e'$ . If in this way we reduce *all* the non-prohibited redexes contained in  $e$ , we call the resulting term  $e^*$  (or  $[e]^*$ ),

which we call the *complete reduction* of  $e$ . (Again, the order in which we perform  $\beta$ -reductions associated with disjoint redexes does not matter.) By definition, we then have  $e \Rightarrow_{\beta} e^*$ .

We now have a similar sequence of lemmas to that contained in the proof of Theorem 1.13. Recall that  $s \rightarrow_{\beta\eta} t$  means  $s \rightarrow_{\beta} t$  or  $s \rightarrow_{\eta} t$ .

**Lemma 1:** For any  $\lambda$ -terms  $s$  and  $t$ ,  $s \rightarrow_{\beta\eta} t$  entails  $s \Rightarrow_{\beta\eta} t$ .

**Proof:** The reduction  $s \rightarrow_{\beta\eta} t$  corresponds to the case in which a single redex  $r$  in the redex trees of  $s$  is reduced, with the result being  $t$ . If the redex  $r$  is not prohibited, then it is immediate that  $s \Rightarrow_{\beta\eta} t$ .

Suppose then that  $r$  is prohibited. Then  $r$  has one of the two forms **(a)** or **(b)** defined earlier. In case **(a)** it has the form of the  $\beta$ -redex  $\lambda x(E(x))(N)$  (with  $x$  not free in  $E$ ), which  $\beta$ -reduces to  $E(N)$ . Note however that if in this term the  $\eta$ -redex  $\lambda x(E(x))$  (which is a proper subterm of  $r$ ) is reduced, the result is also  $E(N)$ .

In case **(b)**, a similar phenomenon occurs. In this case  $r$  has the form of the  $\eta$ -redex  $\lambda z(\lambda x(M)(z))$ , where  $z$  is not free in  $\lambda x(M)$ . This  $\eta$ -reduces to  $\lambda x(M)$ . However, if in this term the  $\beta$ -redex  $\lambda x(M)(z)$  (which is a proper subterm of  $r$ ) is reduced, the result is  $\lambda z(M[z/x])$ . Note that if  $z$  is not free in  $\lambda x(M)$ , then the terms  $\lambda x(M)$  and  $\lambda z(M[z/x])$  are  $\alpha$ -equivalent, and thus the same.

As a consequence, we have

$$\text{If } s \rightarrow_{\beta\eta} t \text{ via the reduction of a prohibited redex } r, \text{ then } s \rightarrow_{\beta\eta} t \text{ via} \quad (23)$$

$$\text{the reduction of a redex } r' \text{ properly contained in } r.$$

Of course, the redex  $r'$  given by (23) may also be prohibited. But then using the principle (23) again, we can find a redex  $r''$  properly contained in  $r'$  such  $s \rightarrow_{\beta\eta} t$  via the reduction of the redex  $r''$ . Because we cannot have an infinite sequence  $r, r', r'', \dots$  of redexes, each a subredex of the previous one on the list, there must be a non-prohibited redex  $r^*$  in  $s$  such that  $s \rightarrow_{\beta\eta} t$  via the reduction of  $r^*$ . Thus,  $s \Rightarrow_{\beta\eta} t$ .  $\square$

**Lemma 2:** For any  $\lambda$ -terms  $s$  and  $t$ ,  $s \Rightarrow_{\beta\eta} t$  entails  $s \rightarrow_{\beta\eta} t$ .

**Proof:** Identical to that given in the proof of Theorem 1.13  $\square$

**Lemma 3:** For any  $\lambda$ -terms  $s, s', t, t'$ , if  $s \Rightarrow_{\beta\eta} s'$  and  $t \Rightarrow_{\beta\eta} t'$ , then  $s[t/v] \Rightarrow_{\beta\eta} s'[t'/v]$ .

**Proof:** Exactly as in Theorem 1.13, but using  $\beta\eta$ -reduction instead of  $\beta$ -reduction.  $\square$

**Lemma 4:** For any  $\lambda$ -terms  $s$  and  $t$ , if  $s \Rightarrow_{\beta\eta} t$  then  $t \Rightarrow_{\beta\eta} s^*$ .

**Proof:** In this proof, in order to deal with both types of reduction at once we will sometimes rewrite the  $\beta$ -redex  $\lambda x(M)(N)$  more abstractly as  $W(M, N)$  and the  $\eta$ -redex  $\lambda x(E(x))$  more abstractly as  $W(E)$ . In this way, all redexes have the form  $W(X, Y, \dots)$  for some  $X, Y, \dots$ . Given a redex  $W(X, Y, \dots)$ , we will write  $W^{\mathbf{r}}(X, Y, \dots)$  for the reduction of this redex. So for example, if  $W(M, N)$  is the  $\beta$ -redex  $\lambda x(M)(N)$ , then  $W^{\mathbf{r}}(M, N)$  is the term  $M[N/x]$ , and if  $W(E)$  is the  $\eta$ -redex  $\lambda x(E(x))$  (with  $x$  not free in  $E$ ), then  $W^{\mathbf{r}}(E)$  is the term  $E$ .

Using this notation, we have the following

- (P1)** If a redex  $r = W(X, Y, \dots)$  is not prohibited and another redex  $r'$  is properly contained in  $r$ , then  $r'$  is properly contained in one of  $X, Y, \dots$
- (P2)** For any redex  $W(X, Y, \dots)$ , if  $X \Rightarrow_{\beta\eta} X', Y \Rightarrow_{\beta\eta} Y', \dots$  then

$$W^{\mathbf{r}}(X, Y, \dots) \Rightarrow_{\beta\eta} W^{\mathbf{r}}(X', Y', \dots).$$

The property **(P1)** follows from our characterization of prohibited redexes. For a  $\beta$ -redex  $\lambda x(M)(N)$ , property **(P2)** amounts to the claim that if  $M \Rightarrow_{\beta\eta} M'$  and  $N \Rightarrow_{\beta\eta} N'$ , then  $M[N/x] \Rightarrow_{\beta\eta} M'[N'/x]$ , which follows from Lemma 3. For an  $\eta$ -redex  $\lambda x(E(x))$ , property **(P2)** amounts to the trivial claim that if  $E \Rightarrow_{\beta\eta} E'$  then  $E \Rightarrow_{\beta\eta} E'$ . So both **(P1)** and **(P2)** are true in general. They will be useful in what follows.

As in Theorem 1.13, Lemma 4 is proved by induction on  $\text{height}(s)$ , where  $\text{height}(s)$  is the length of the longest branch in the redex trees of  $s$  (this includes prohibited redexes.) As before, the base case is trivial.

Suppose now that Lemma 4 holds for terms  $s'$  with  $\text{height}(s') \leq l$ , and consider a term  $s$  with  $\text{height}(s) = l + 1$ . Suppose initially that  $s$  is a redex, i.e.,

$$s = W(X, Y, \dots).$$

Then note that we have

$$s^* = W^{\mathbf{r}}(X^*, Y^*, \dots).$$

Suppose next that  $s \Rightarrow_{\beta} t$ . There are two cases to consider

- (Case 1)** In the reduction  $s \Rightarrow_{\beta} t$ , the redex  $s$  itself is among the redexes that are reduced.
- (Case 2)** In the reduction  $s \Rightarrow_{\beta} t$ , the redex  $s$  itself is not among the redexes that are reduced.

(Note of course that if in the reduction  $s \Rightarrow_{\beta} t$  the redex  $s$  is reduced, it will be the last redex reduced, as it lies at the very bottom of the redex trees for  $s$ .)

Making use of **(P1)**, in Case 1 the term  $t$  has the form  $W^{\mathbf{r}}(X', Y', \dots)$  for some  $X \Rightarrow_{\beta\eta} X', Y \Rightarrow_{\beta\eta} Y', \dots$ , while in Case 2 the term  $t$  has the form  $W(X', Y', \dots)$  for some  $X \Rightarrow_{\beta\eta} X', Y \Rightarrow_{\beta\eta} Y', \dots$ . In either case,  $\text{height}(X), \text{height}(Y), \dots$  are all strictly less than  $\text{height}(s)$ , and so the inductive hypothesis holds of  $X, Y, \dots$ . Thus we have  $X' \Rightarrow_{\beta\eta} X^*, Y' \Rightarrow_{\beta\eta} Y^*, \dots$

In Case 1 we then have  $t = W^{\mathbf{r}}(X', Y', \dots)$  and  $X' \Rightarrow_{\beta\eta} X^*, Y' \Rightarrow_{\beta\eta} Y^*, \dots$ , and so using **(P2)** we have  $t \Rightarrow_{\beta\eta} W^{\mathbf{r}}(X^*, Y^*, \dots) = s^*$ .

In Case 2 we have  $t = W(X', Y', \dots)$ ,  $X' \Rightarrow_{\beta\eta} X^*, Y' \Rightarrow_{\beta\eta} Y^*, \dots$  and so  $t \Rightarrow_{\beta\eta} W^{\mathbf{r}}(X^*, Y^*) = s^*$ . (The reduction  $t \Rightarrow_{\beta\eta} W^{\mathbf{r}}(X^*, Y^*) = s^*$  is obtained by starting with  $t = W(X', Y', \dots)$ , reducing  $X' \Rightarrow_{\beta\eta} X^*, Y' \Rightarrow_{\beta\eta} Y^*, \dots$  to obtain  $W(X^*, Y^*, \dots)$ , then finally reducing the redex  $W$  to obtain  $W^{\mathbf{r}}(X^*, Y^*, \dots) = s^*$ .)

Therefore, in either case we have  $t \Rightarrow_{\beta\eta} s^*$  as desired. To complete the induction, we must finally consider the case in which  $s$  is not a redex. This is dealt with as in Theorem 1.13.  $\square$



Given Lemma 4, the proof of Lemma 5 and remainder of the proof is the same as before, with only trivial notational changes (e.g., replacing  $\beta$  wherever it appears with  $\beta\eta$ .)  $\square$

**Theorem 1.31: Strong Normalizability of  $\lambda_0^{\text{type}}$  under  $\beta$ -reduction.**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$$

**Proof of Theorem 1.31.**

In this proof, all  $\lambda$ -terms are understood to be terms of the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ .

For each fundamental type  $\tau$ , we add a new atomic, closed term  $d^{\tau}$  of type  $\tau$  to the typed  $\lambda$ -calculus. For any compound type  $\tau_1 \rightarrow \tau_2$ , we let  $d^{\tau_1 \rightarrow \tau_2}$  be an abbreviation for  $\lambda x^{\tau_1} (d^{\tau_2})$ . In this way,  $d^{\tau}$  is inductively defined for all simple types  $\tau$ , and is always a closed term of type  $\tau$ , already in  $\beta$ -normal form. (Without expanding the  $\lambda$ -calculus in this way, it is not true that for every type  $\tau$  there is a closed term of type  $\tau$ .) Showing that strong normalization holds for this expanded typed  $\lambda$ -calculus obviously entails that it holds for our original typed  $\lambda$ -calculus.

**Definition:** A  $\lambda$ -term  $e$  is *strongly normalizing* if there is no infinite sequence  $e_0, e_1, e_2, \dots$  of  $\lambda$ -terms with  $e = e_0$  such that each  $e_{i+1}$  may be obtained from  $e_i$  by a single step  $\beta$ -reduction.

For all types  $\tau$ , we define a set  $\llbracket \tau \rrbracket$  of closed  $\lambda$ -terms of type  $\tau$  as follows

**Definition:** For every type  $\tau$ ,  $\llbracket \tau \rrbracket$  is given as follows

- (i) If  $\tau$  is fundamental,  $\llbracket \tau \rrbracket$  is the set of strongly normalizing, closed terms of type  $\tau$ .
- (ii) If  $\tau$  has the form  $\tau_1 \rightarrow \tau_2$ ,  $\llbracket \tau \rrbracket$  is the set of closed terms  $e$  of type  $\tau$  such that for all  $e' \in \llbracket \tau_1 \rrbracket$ , we have that  $e(e') \in \llbracket \tau_2 \rrbracket$ .

It is easily seen (by induction on the construction of types) that every type can be written uniquely in the form

$$\tau_1 \rightarrow (\tau_2 \rightarrow (\tau_3 \rightarrow \dots \rightarrow (\tau_n \rightarrow \sigma) \dots))$$

where  $\sigma$  is a fundamental type. We use this fact freely in what follows. The following alternate characterization of  $\llbracket \tau \rrbracket$  is useful

**Lemma 1:** If a closed term  $e$  has type  $\tau = \tau_1 \rightarrow (\tau_2 \rightarrow (\tau_3 \rightarrow \dots \rightarrow (\tau_n \rightarrow \sigma) \dots))$  where  $\sigma$  is fundamental, then  $e$  is in  $\llbracket \tau \rrbracket$  iff for all closed terms  $e_1 \in \llbracket \tau_1 \rrbracket, e_2 \in \llbracket \tau_2 \rrbracket, \dots, e_n \in \llbracket \tau_n \rrbracket$ , we have that  $e(e_1) \dots (e_n)$  is strongly normalizing.

**Proof:** By induction on  $\tau$ . If  $\tau$  is fundamental, the result is immediate. Suppose then that  $\tau$  has the form  $\sigma_1 \rightarrow \sigma_2$  and the lemma holds of  $\sigma_2$ . Suppose further that  $\sigma_2$  has the form

$\tau_1 \rightarrow (\dots \rightarrow (\tau_n \rightarrow \rho)\dots)$  where  $\rho$  is fundamental, so that  $\sigma_1 \rightarrow \sigma_2$  can be written in the form  $\sigma_1 \rightarrow (\tau_1 \rightarrow (\dots \rightarrow (\tau_n \rightarrow \rho)\dots))$ . Then for any closed term  $e$  of type  $\sigma_1 \rightarrow \sigma_2$ , we have

$$\begin{aligned} e \in \llbracket \sigma_1 \rightarrow \sigma_2 \rrbracket & \text{ iff for all closed terms } e' \in \llbracket \sigma_1 \rrbracket, e(e') \in \llbracket \sigma_2 \rrbracket \text{ (by definition)} \\ & \text{ iff for all closed terms } e' \in \llbracket \sigma_1 \rrbracket, e_1 \in \llbracket \tau_1 \rrbracket, \dots, e_n \in \llbracket \tau_n \rrbracket, e(e')(e_1)\dots(e_n) \\ & \text{ is strongly normalizing (applying inductive hypothesis to } e(e')). \end{aligned}$$

as desired. This completes the induction.  $\square$ .

The following technical lemma forms the heart of the proof.

**Lemma 2:** For any type  $\tau$ ,

- (i) If  $t \in \llbracket \tau \rrbracket$  and  $t$  is strongly normalizing, and  $e$  is a closed term with  $e \in \llbracket \sigma \rrbracket$ , then  $\lambda x^\tau(e)(t) \in \llbracket \sigma \rrbracket$ .
- (ii) If  $t \in \llbracket \tau \rrbracket$ , and  $e$  is term of type  $\sigma$  with one free variable  $x^\tau$  and  $e[t/x^\tau] \in \llbracket \sigma \rrbracket$ , then  $\lambda x^\tau(e)(t) \in \llbracket \sigma \rrbracket$ .

**Proof:** We first prove (i). Suppose the type  $\sigma$  can be written as

$$\sigma_1 \rightarrow (\sigma_2 \rightarrow (\sigma_3 \rightarrow \dots \rightarrow (\sigma \rightarrow \rho)\dots)),$$

where  $\rho$  is a fundamental type. Suppose then that  $t \in \llbracket \tau \rrbracket$ ,  $t$  is strongly normalizing,  $e \in \llbracket \sigma \rrbracket$ , and  $\lambda x^\tau(e)(t) \notin \llbracket \sigma \rrbracket$ . Then by Lemma 1, for some terms  $e_1 \in \llbracket \tau_1 \rrbracket, e_2 \in \llbracket \tau_2 \rrbracket, \dots, e_n \in \llbracket \tau_n \rrbracket$  we have that

$$\lambda x^\tau(e)(t)(e_1)\dots(e_n) \tag{24}$$

is not strongly normalizing, so that we may perform an infinite sequence of  $\beta$ -reductions starting with (24). Fix some infinite sequence of  $\beta$ -reductions that begins with the term (24).

It is easily seen that the only redexes contained in a term of the form (24) are the redex  $\lambda x^\tau(e)(t)$ , and any redexes that occur as (not necessarily proper) subexpressions of  $e, t, e_1, \dots, e_n$ . Thus, the first  $\beta$ -reduction must be either performed on (a) the initial redex  $\lambda x^\tau(e)(t)$ , or (b) a (not necessarily proper) subterm of one of  $e, t, e_1, \dots, e_n$ . If the first reduction is not of form (a), then the newly reduced term can be written in the form  $\lambda x^\tau(e^1)(t^1)(e_1^1)\dots(e_n^1)$  (where all but one of  $e^1, t^1, e_1^1, \dots, e_n^1$  is identical with the corresponding  $e, t, e_1, \dots, e_n$ ), and so the next  $\beta$ -reduction must be either performed on (a) the initial subterm  $\lambda x^\tau(e^1)(t^1)$ , or (b) a (not necessarily proper) subterm of one of  $e^1, t^1, e_1^1, \dots, e_n^1$ . Iterating this reasoning, we see that in any infinite reduction sequence, either (c) a  $\beta$ -reduction involving the initial subterm  $\lambda x^\tau(e^*)(t^*)$  with  $e \rightarrow_\beta e^*$  and  $t \rightarrow_\beta t^*$  is eventually performed at some point, or (d) the  $i$ th term of the reduction sequence has the form  $\lambda x^\tau(e^i)(t^i)(e_1^i)\dots(e_n^i)$ , where letting  $e^0 = e, t^0 = t, e_1^0 = e_1, \dots, e_n^0 = e_n$ , we have that for each  $i > 0$ , exactly one of the terms  $e^i, t^i, e_1^i, \dots, e_n^i$  is obtained by a one-step  $\beta$ -reduction from  $e^{i-1}, t^{i-1}, e_1^{i-1}, \dots, e_n^{i-1}$ , and the remainder of the  $e^i, t^i, e_1^i, \dots, e_n^i$  are identical to the corresponding  $e^{i-1}, t^{i-1}, e_1^{i-1}, \dots, e_n^{i-1}$ . We consider these two cases (c) and (d) in turn.

In case (c), because the term  $e$  is closed, so too the term  $e^*$  is closed, and the redex  $\lambda x^\tau(e^*)(t^*)$  thus  $\beta$ -reduces to  $e^*$ . Thus a term of the form

$$e^*(e_1^*)\dots(e_n^*)$$

with  $e \rightarrow_\beta e^*$ ,  $e_1 \rightarrow_\beta e_1^*$ , ...,  $e_n \rightarrow_\beta e_n^*$  occurs in the infinite  $\beta$ -reduction sequence that starts with (24). Thus  $e(e_1)\dots(e_n)$  is also not strongly normalizing, as  $e(e_1)\dots(e_n) \rightarrow_\beta e^*(e_1^*)\dots(e_n^*)$ . But if  $e(e_1)\dots(e_n)$  is not strongly normalizing, then because  $e_1 \in \llbracket \tau_1 \rrbracket$ ,  $e_2 \in \llbracket \tau_2 \rrbracket$ , ...,  $e_n \in \llbracket \tau_n \rrbracket$ , it follows from Lemma 1 that  $e \notin \llbracket \sigma \rrbracket$ , which is a contradiction.

Consider then case (d). In this case, our infinite reduction sequence involves nothing more than reductions of the subexpressions  $e, t, e_1, \dots, e_n$  individually. Thus, at least one of the terms  $e, t, e_1, \dots, e_n$ , must fail to be strongly normalizing. By assumption,  $t$  is strongly normalizing, and thus one of  $e, e_1, \dots, e_n$  must fail to be strongly normalizing. Thus  $e(e_1)\dots(e_n)$  is not strongly normalizing, which as in case (c) contradicts  $e \in \llbracket \sigma \rrbracket$ . This completes the proof of (i).

The proof of (ii) is similar. This time, suppose that  $t \in \llbracket \tau \rrbracket$ ,  $e[t/x^\tau] \in \llbracket \sigma \rrbracket$ , and  $\lambda x^\tau(e)(t) \notin \llbracket \sigma \rrbracket$ . Then for some terms  $e_1 \in \llbracket \tau_1 \rrbracket$ ,  $e_2 \in \llbracket \tau_2 \rrbracket$ , ...,  $e_n \in \llbracket \tau_n \rrbracket$ , we have that

$$\lambda x^\tau(e)(t)(e_1)\dots(e_n) \tag{25}$$

is not strongly normalizing. Fix some infinite sequence of  $\beta$ -reductions that begin with the term (25). As before, either (c) a  $\beta$ -reduction involving the initial subterm  $\lambda x^\tau(e^*)(t^*)$  with  $e \rightarrow_\beta e^*$  and  $t \rightarrow_\beta t^*$  is eventually performed at some point, or (d) the infinite reduction sequence involves nothing more than reductions of the subexpressions  $e, t, e_1, \dots, e_n$  individually.

In case (c), the redex  $\lambda x^\tau(e^*)(t^*)$   $\beta$ -reduces to  $e^*[t^*/x^\tau]$ , and thus a term of the form

$$e^*[t^*/x^\tau](e_1^*)\dots(e_n^*)$$

with  $e \rightarrow_\beta e^*$ ,  $t \rightarrow_\beta t^*$ ,  $e_1 \rightarrow_\beta e_1^*$ , ...,  $e_n \rightarrow_\beta e_n^*$  occurs in the infinite  $\beta$ -reduction sequence that starts with (25). Thus  $e[t/x^\tau](e_1)\dots(e_n)$  is not strongly normalizing (as  $e[t/x^\tau](e_1)\dots(e_n) \rightarrow_\beta e^*[t^*/x^\tau](e_1^*)\dots(e_n^*)$ ), contradicting  $e[t/x^\tau] \in \llbracket \sigma \rrbracket$ .

In case (d), our infinite reduction sequence involves nothing more than reductions of the subexpressions  $e, t, e_1, \dots, e_n$  individually. Thus, at least one of the terms  $e, t, e_1, \dots, e_n$  must fail to be strongly normalizing. In any such case, the term  $e[t/x^\tau](e_1)\dots(e_n)$  also fails to be strongly normalizing so long as we suppose, as we have, that the variable  $x^\tau$  actually occurs in  $e$ , so that  $t$  actually occurs as a subexpression of  $e[t/x^\tau]$ . (We also invoke the easily proved fact that if there is an infinite reduction sequence starting from  $e$ , then there is an infinite reduction sequence starting from  $e[t/x^\tau]$ .) But  $e[t/x^\tau](e_1)\dots(e_n)$  not being strongly normalizing contradicts  $e[t/x^\tau] \in \llbracket \sigma \rrbracket$ .  $\square$

**Lemma 3:** For any type  $\tau$ ,

- (i)  $d^\tau \in \llbracket \tau \rrbracket$
- (ii) for all  $t \in \llbracket \tau \rrbracket$ ,  $t$  is strongly normalizing.

**Proof:** By induction on  $\tau$ . If  $\tau$  is fundamental, (i) and (ii) are immediate from the fact that  $d^\tau$  is in  $\beta$ -normal form and the definition of  $\llbracket \tau \rrbracket$  for fundamental  $\tau$ .

Suppose by inductive hypothesis that (i) and (ii) are true for types  $\tau_1$  and  $\tau_2$ , and consider the type  $\tau = \tau_1 \rightarrow \tau_2$ . Suppose that  $\tau_2$  may be written  $\sigma_1 \rightarrow (\sigma_2 \rightarrow \dots \rightarrow (\sigma_n \rightarrow \rho))$ , where  $\rho$  is fundamental. We may also suppose by inductive hypothesis that (i) and (ii) hold of  $\sigma_1, \dots, \sigma_n$ .

To prove (i) for  $\tau$ , note that by definition  $d^\tau \in \llbracket \tau \rrbracket$  iff for all  $t \in \llbracket \tau_1 \rrbracket$ , we have  $d^\tau(t) \in \llbracket \tau_2 \rrbracket$ . Because  $d^\tau$  is  $\lambda x^{\tau_1}(d^{\tau_2})$ , we have that  $d^\tau \in \llbracket \tau \rrbracket$  iff

$$\text{for all } t \in \llbracket \tau_1 \rrbracket, \lambda x^{\tau_1}(d^{\tau_2})(t) \in \llbracket \tau_2 \rrbracket. \quad (26)$$

Because by inductive hypothesis we know that all elements of  $\llbracket \tau_1 \rrbracket$  are strongly normalizing and that  $d^{\tau_2} \in \llbracket \tau_2 \rrbracket$ , that (26) is true follows from Lemma 2(i).

To prove (ii) for  $\tau$ , note that by inductive hypothesis we have  $d^{\tau_1} \in \llbracket \tau_1 \rrbracket, d^{\sigma_1} \in \llbracket \sigma_1 \rrbracket, \dots, d^{\sigma_n} \in \llbracket \sigma_n \rrbracket$ . Thus, from  $t \in \llbracket \tau \rrbracket$  and Lemma 1, we have that  $t(d^{\tau_1})(d^{\sigma_1})\dots(d^{\sigma_n})$  is strongly normalizing. Because any subterm of a strongly normalizing term is strongly normalizing, it follows that  $t$  is strongly normalizing.  $\square$

**Lemma 4:** If  $e$  is a (not necessarily closed)  $\lambda$ -term of type  $\sigma$  with free variables  $x_1^{\sigma_1}, x_2^{\sigma_2}, \dots, x_n^{\sigma_n}$  and  $t_1, \dots, t_n$  are closed  $\lambda$ -terms of type  $\sigma_1, \dots, \sigma_n$  respectively with  $t_1 \in \llbracket \sigma_1 \rrbracket, \dots, t_n \in \llbracket \sigma_n \rrbracket$ , then

$$e[t_1/x_1^{\sigma_1}]\dots[t_n/x_n^{\sigma_n}] \in \llbracket \sigma \rrbracket.$$

(Note that because the terms  $t_i$  are closed, the order of the substitutions  $t_1/x_1^{\sigma_1}, \dots, t_n/x_n^{\sigma_n}$  does not matter.)

**Proof:** By induction on the construction of  $e$ .

For the first base case, suppose  $e$  is just the variable  $x^\sigma$ . We must prove that if  $t$  is a closed  $\lambda$ -term of type  $\sigma$  with  $t \in \llbracket \sigma \rrbracket$ , then  $x[t/x] \in \llbracket \sigma \rrbracket$ , i.e.,  $t \in \llbracket \sigma \rrbracket$ , which is immediate.

For the second base case, suppose  $e$  is  $d^\sigma$  for some fundamental type  $\sigma$ . Then  $e$  has no free variables, and it suffices to show  $d^\sigma \in \llbracket \sigma \rrbracket$ , which we have from Lemma 3.

Next suppose that  $e$  has the form  $e_1(e_2)$  and has the type  $\sigma$ , so that  $e_1$  has the type  $\sigma \rightarrow \sigma'$  and  $e_2$  has the type  $\sigma'$  for some type  $\sigma'$ . Suppose also that  $e$  has free variables  $x_1^{\sigma_1}, x_2^{\sigma_2}, \dots, x_n^{\sigma_n}$  and that  $t_1, \dots, t_n$  are closed  $\lambda$ -terms of type  $\sigma_1, \dots, \sigma_n$  with  $t_1 \in \llbracket \sigma_1 \rrbracket, \dots, t_n \in \llbracket \sigma_n \rrbracket$ . Then by inductive hypothesis we have  $e_1[t_1/x_1]\dots[t_n/x_n] \in \llbracket \sigma' \rightarrow \sigma \rrbracket$  and  $e_2[t_1/x_1]\dots[t_n/x_n] \in \llbracket \sigma' \rrbracket$ . From this and the definition of  $\llbracket \sigma' \rightarrow \sigma \rrbracket$ , it follows that

$$e_1[t_1/x_1]\dots[t_n/x_n](e_2[t_1/x_1]\dots[t_n/x_n]) \in \llbracket \sigma \rrbracket.$$

From the basic properties of substitution we then have

$$e_1(e_2)[t_1/x_1]\dots[t_n/x_n] \in \llbracket \sigma \rrbracket$$

as desired.

Finally, suppose that  $e$  has the form  $\lambda z^\sigma(u)$  where  $u$  has the type  $\sigma'$ , so that  $e$  has the type  $\sigma \rightarrow \sigma'$ . Suppose that  $e = \lambda z^\sigma(u)$  has free variables  $x_1^{\sigma_1}, x_2^{\sigma_2}, \dots, x_n^{\sigma_n}$  (all of course distinct from  $z^\sigma$ ) and that  $t_1 \in \llbracket \sigma_1 \rrbracket, \dots, t_n \in \llbracket \sigma_n \rrbracket$ . To start, assume that the variable  $z^\sigma$  explicitly occurs free in  $u$ . The free variables of  $u$  are then  $x_1^{\sigma_1}, x_2^{\sigma_2}, \dots, x_n^{\sigma_n}$  and  $z^\sigma$ . By inductive hypothesis we have

$$\text{for all } t \in \llbracket \sigma \rrbracket, \quad u[t_1/x_1^{\sigma_1}] \dots [t_n/x_n^{\sigma_n}][t/z^\sigma] \in \llbracket \sigma' \rrbracket.$$

Applying Lemma 2(ii) (with  $u[t_1/x_1^{\sigma_1}] \dots [t_n/x_n^{\sigma_n}]$  playing the role of  $e$ ), we have

$$\text{for all } t \in \llbracket \sigma \rrbracket, \quad \lambda z^\sigma(u[t_1/x_1] \dots [t_n/x_n])(t) \in \llbracket \sigma' \rrbracket$$

from which it follows from the definition of  $\llbracket \sigma \rightarrow \sigma' \rrbracket$  that

$$\lambda z^\sigma(u[t_1/x_1] \dots [t_n/x_n]) \in \llbracket \sigma \rightarrow \sigma' \rrbracket$$

and because the terms  $t_i$  are closed,

$$\lambda z^\sigma(u)[t_1/x_1] \dots [t_n/x_n] \in \llbracket \sigma \rightarrow \sigma' \rrbracket$$

as desired.

It suffices to consider the case in which  $e$  has the form  $\lambda z^\sigma(u)$  and the variable  $z^\sigma$  does *not* occur free in  $u$ . In this case (using the same notation), by inductive hypothesis we have

$$u[t_1/x_1] \dots [t_n/x_n] \in \llbracket \sigma' \rrbracket.$$

Using Lemma 2(i) (and invoking the fact proved in Lemma 3 that all the elements of  $\llbracket \sigma \rrbracket$  are strongly normalizing), we have as before that

$$\text{for all } t \in \llbracket \sigma \rrbracket, \quad \lambda z^\sigma(u[t_1/x_1] \dots [t_n/x_n])(t) \in \llbracket \sigma' \rrbracket.$$

from which it again follows that

$$\lambda z^\sigma(u)[t_1/x_1] \dots [t_n/x_n] \in \llbracket \sigma \rightarrow \sigma' \rrbracket$$

as desired.  $\square$

It follows from Lemma 4 that if  $e$  is a *closed*  $\lambda$ -term of type  $\sigma$ , then  $e \in \llbracket \sigma \rrbracket$ . From Lemma 3 it then follows that  $e$  is strongly normalizing. So all closed  $\lambda$ -terms are strongly normalizing. Let  $e$  be an open  $\lambda$ -term with free variables  $x_1^{\sigma_1}, x_2^{\sigma_2}, \dots, x_n^{\sigma_n}$ . Then the  $\lambda$ -term  $\lambda x_1^{\sigma_1} x_2^{\sigma_2} \dots x_n^{\sigma_n}(e)$  is closed and therefore strongly normalizing. But because every subterm of a strongly normalizing term is strongly normalizing, it follows that  $e$  is strongly normalizing. So all  $\lambda$ -terms are strongly normalizing.  $\square$

**Theorem 1.32: The Church-Rosser Theorem for  $\beta$ -reductions in  $\lambda_0^{\text{type}}$ .**

In the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$ , if  $X \rightarrow_{\beta} Y_1$  and  $X \rightarrow_{\beta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ .

**Theorem 1.34: The Church-Rosser Theorem for  $\beta\eta$ -reductions in  $\lambda_0^{\text{type}}$ .**

In the typed  $\lambda$ -calculus  $\lambda_0$ , if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

**Proof of Theorems 1.32 and 1.34.**

To prove Theorems 1.32 and 1.34, it suffices to use the proofs of Theorems 1.13 and 1.20 and note that nothing more than notational alterations are required in the typed case. In particular, every term used in these proofs can be typed appropriately if required.  $\square$

**Theorem 1.43: Soundness and Completeness of  $\text{TR}_0$ .**

Suppose  $E$  is an untyped  $\lambda_0$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) It is possible to assign simple types (i.e., types of  $\lambda_0^{\text{type}}$ ) to all the variables of  $E$  in such a way that the free variables  $v_1, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \dots, \tau_n$ , and  $E$  itself has type  $X$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\text{TR}_0$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

**Proof of Theorem 1.43.**

First we prove (i)  $\rightarrow$  (ii). The proof is by induction on the construction of the  $\lambda$ -term  $E$ . Let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ . We assume without loss of generality that  $E$  obeys the Barendregt variable convention. For suppose  $E$  is a term that does not obey the Barendregt variable convention, and let  $E'$  be a term  $\alpha$ -equivalent to  $E$  that does obey the Barendregt variable convention. Because  $E$  and  $E'$  are  $\alpha$ -equivalent, if  $E$  satisfies (i), then so does  $E'$ . Thus, if (i)  $\rightarrow$  (ii) holds for  $E'$ , we have a proof of  $\Gamma \vdash E' : X$  for some  $\Gamma$  with  $\bar{\Gamma} \subseteq \Gamma$ . Silent replacement of  $\alpha$ -equivalent terms in the final step of this proof then gives  $\Gamma \vdash E : X$ , and so

**(ii)** holds, as desired.

For the base case of the induction, suppose that  $E$  is a variable  $v$ . Then  $E$  can be assigned any type  $X$  whatsoever. Because there is a proof of  $v : X \vdash v : X$  for any  $X$ , **(i)**  $\rightarrow$  **(ii)** holds in this case.

Suppose next that  $E$  has the form  $E_1(E_2)$ . Let us suppose that consistent with the variable declarations in  $\bar{\Gamma}$ , it is possible to type the remaining variables of  $E$  so that  $E$  is assigned the type  $X$ . In this situation, we must have that for some type  $Y$ ,  $E_1$  is assigned type  $Y \rightarrow X$  and  $E_2$  is assigned type  $Y$ . Then by inductive hypothesis we have

$$\Gamma_1 \vdash E_1 : Y \rightarrow X \quad \text{and} \quad \Gamma_2 \vdash E_2 : Y \quad (27)$$

for contexts  $\Gamma_1$  and  $\Gamma_2$  such that for every free variable  $v_i$  of  $E_1$ , the context  $\Gamma_1$  contains  $v_i : \tau_i$ , and for every free variable  $v_i$  of  $E_2$ , the context  $\Gamma_2$  contains  $v_i : \tau_i$ . By Theorem 1.41, we may in fact assume that  $\Gamma_1$  and  $\Gamma_2$  consist *only* of typing declarations for the free variables of  $E_1$  and  $E_2$  respectively (and thus  $\Gamma_1, \Gamma_2 \subseteq \bar{\Gamma}$ ), and that the only other variable declarations that appear anywhere else in the typing derivations (27) are typing declarations for bound variables of  $E_1$  and  $E_2$ . Because  $E_1$  and  $E_2$  obey the Barendregt variable convention, the bound variables of  $E_1$  and  $E_2$  are distinct from the free variables of  $E_1(E_2)$ , and so these additional variable declarations all involve variables distinct from  $v_1, \dots, v_n$ .

By adding additional (redundant) variable declarations to the contexts in the typing derivations (27), we may then produce typing derivations

$$\bar{\Gamma} \vdash E_1 : Y \rightarrow X \quad \text{and} \quad \bar{\Gamma} \vdash E_2 : Y. \quad (28)$$

(Note that we must invoke the fact that any additional variable declarations in (27) that occur above the final conclusion involve variables other than  $v_1, \dots, v_n$ . If this were not true there would be a chance that adding redundant variable declarations to the typing constructions in (27) would produce a context with a variable declared in two different ways, which would be ungrammatical. By working only with terms that obey the Barendregt variable convention we avoid this possibility.)

From (28), we may infer  $\bar{\Gamma} \vdash E : X$  using the (App) rule. So **(i)**  $\rightarrow$  **(ii)** holds in this case.

Suppose finally that  $E$  has the form  $\lambda v(E')$  and that the free variables of  $\lambda v(E')$  are  $v_1, \dots, v_n$ , so that the free variables of  $E'$  are either  $v, v_1, \dots, v_n$  or  $v_1, \dots, v_n$ , depending on whether  $v$  explicitly appears free in  $E'$ . Suppose also that consistent with the variable declarations of  $\bar{\Gamma}$ , the remaining variables of  $\lambda v(E')$  can be assigned types in such a way that  $\lambda v(E')$  is assigned type  $X$ . In this situation, for some types  $Y$  and  $Z$  the type  $X$  will have the form  $Y \rightarrow Z$ , the term  $E'$  will have the type  $Z$ , and the variable  $v$  will have the type  $Y$ . Then by inductive hypothesis and Theorem 1.41 there is a typing derivation with conclusion

$$\bar{\Gamma}, v : Y \vdash E' : Z \quad \text{or} \quad \bar{\Gamma} \vdash E' : Z \quad (29)$$

depending on whether  $v$  explicitly appears free in  $E'$  or not. In the second case of (29), the variable  $v$  does not occur bound in  $E'$  and so by Theorem 1.41 we may assume that  $v$  is not declared anywhere in the typing derivation  $\bar{\Gamma} \vdash E' : Z$ . Thus we may freely add  $v : Y$  as a (redundant) variable declaration throughout the proof of  $\bar{\Gamma} \vdash E' : Z$ , obtaining a typing derivation with conclusion  $\bar{\Gamma}, v : Y \vdash E' : Z$ ; i.e., the first case of (29). So either way we have a typing derivation with conclusion  $\bar{\Gamma}, v : Y \vdash E' : Z$ . Applying the (Abs) rule then gives

$$\bar{\Gamma} \vdash \lambda v(E') : Y \rightarrow Z$$

as desired, and **(i)**  $\rightarrow$  **(ii)** holds in this case. This completes the proof that **(i)**  $\rightarrow$  **(ii)** holds in general.

Now we prove **(ii)**  $\rightarrow$  **(i)**. The proof is by induction on the construction of typing derivations. (We do not need to assume anymore that any particular terms obey the Barendregt variable convention.) If  $\Gamma \vdash E : X$  has the form of an initial sequent  $\Delta, v : X \vdash v : X$ , then it suffices to show that it is consistent with  $v$  being assigned the type  $X$  that  $v$  is assigned the type  $X$ , which is trivial. Thus **(ii)**  $\rightarrow$  **(i)** holds in this case.

Suppose next that  $\Gamma \vdash E_1(E_2) : X$  is the result of an (App) inference from

$$\Gamma \vdash E_1 : Y \rightarrow X \quad \text{and} \quad \Gamma \vdash E_2 : Y$$

Let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, \dots, v_n : \tau_n$ , let  $\bar{\Gamma}_1$  be the subset of variable declarations in  $\bar{\Gamma}$  that declare types for free variables of  $E_1$ , and  $\bar{\Gamma}_2$  be the subset of variable declarations in  $\bar{\Gamma}$  that declare types for free variables of  $E_2$ . (Every typing declaration in  $\bar{\Gamma}$  will then be in one or both of  $\bar{\Gamma}_1$  and  $\bar{\Gamma}_2$ , as every free variable of  $E_1(E_2)$  occurs either free in  $E_1$  or free in  $E_2$ .)

By inductive hypothesis, it then follows that consistent with the free variables of  $E_1$  being typed in accordance with the variable declarations in  $\bar{\Gamma}_1$ , the term  $E_1$  can be assigned the type  $Y \rightarrow X$ , and consistent with the the free variables of  $E_2$  being typed in accordance with the variable declarations in  $\bar{\Gamma}_2$ , the term  $E_2$  can be assigned the type  $Y \rightarrow X$ . Because both  $\bar{\Gamma}_1$  and  $\bar{\Gamma}_2$  are contained in  $\bar{\Gamma}$ , there is no conflict between the variable declarations of  $\bar{\Gamma}_1$  and  $\bar{\Gamma}_2$ , and it follows that consistent with the free variables of  $E_1(E_2)$  being typed in accordance with the variable declarations in  $\bar{\Gamma}$ , the term  $E_1(E_2)$  can be assigned the type  $X$ , as desired.

Suppose finally that  $\Gamma \vdash \lambda v(E) : X \rightarrow Y$  is the result of an (Abs) inference from

$$\Gamma, v : X \vdash E : Y \tag{30}$$

where if  $v_1, \dots, v_n$  are the free variables of  $\lambda v(E)$ , then  $\Gamma$  includes the typing declarations  $v_1 : \tau_1, \dots, v_n : \tau_n$ . Again, let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, \dots, v_n : \tau_n$ .

If  $v$  occurs free in  $E$ , then the free variables of  $E$  are  $v_1, \dots, v_n, v$ , and so from (30) and inductive hypothesis, we have that consistent with the variable declarations in  $\bar{\Gamma}, v : X$  it is possible to assign  $E$  the type  $Y$ . It follows that consistent with the variable declarations in  $\bar{\Gamma}$  it is possible to assign  $\lambda v(E)$  the type  $X \rightarrow Y$ , as desired. If  $v$  does not occur free in  $E$ ,



then by inductive hypothesis and (30) we have that consistent with the variable declarations in  $\bar{\Gamma}$  it is possible to assign  $E$  the type  $Y$ . Again, this means that consistent with the variable declarations in  $\bar{\Gamma}$  it is possible to assign  $\lambda v(E)$  the type  $X \rightarrow Y$ , as desired. This completes the proof that **(ii)**  $\rightarrow$  **(i)** holds in general.

**Theorem 1.49: Soundness and Completeness of  $\text{TR}_0^{\text{bt}}$ .**

Suppose  $E$  is a bound-typed  $\lambda$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) If the variables  $v_1, v_2, \dots, v_n$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, then  $E$  must be assigned type  $\sigma$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\text{TR}_0^{\text{bt}}$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

**Proof of Theorem 1.49.**

The proof is largely the same as that of Theorem 1.43. The proof that **(i)**  $\rightarrow$  **(ii)** is by induction on the construction of the  $\lambda$ -term  $E$ . Let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ . As in Theorem 1.43, we may assume that  $E$  obeys the Barendregt variable convention.

The base case in which  $E$  is a variable is trivial. Suppose then that  $E$  has the form  $E_1(E_2)$ . Let us suppose that if the free variables in  $E$  are typed in accordance with  $\bar{\Gamma}$ , the term  $E$  must receive type  $X$ . Because the type of  $E$  and all its subterms are uniquely determined by the typing declarations in  $\bar{\Gamma}$ , these typing declarations also uniquely determine the types of  $E_1$  and  $E_2$ . Suppose then that if the free variables in  $E_1$  and  $E_2$  are typed in accordance with  $\bar{\Gamma}$ , that  $E_1$  must be assigned the type  $Y \rightarrow X$  and  $E_2$  must be assigned type  $Y$ . Then by inductive hypothesis we have

$$\Gamma_1 \vdash E_1 : Y \rightarrow X \quad \text{and} \quad \Gamma_2 \vdash E_2 : Y$$

for contexts  $\Gamma_1$  and  $\Gamma_2$  such that for every free variable  $v_i$  of  $E_1$ ,  $\Gamma_1$  contains  $v_i : \tau_i$ , and for every free variable  $v_i$  of  $E_2$ ,  $\Gamma_2$  contains  $v_i : \tau_i$ . Precisely as in the proof of Theorem 1.43 we can then argue that

$$\bar{\Gamma} \vdash E_1(E_2) : X$$

So **(i)**  $\rightarrow$  **(ii)** holds in this case.

Suppose finally that  $E$  has the form  $\lambda v^Y(E')$  for some type  $Y$ , so that  $E'$  has free variables  $v_1, \dots, v_n$  and possibly  $v$ . Assume also that given the typing declarations in  $\bar{\Gamma}$ ,  $\lambda v^Y(E')$  must be assigned the type  $Y \rightarrow Z$ . If  $v$  occurs free in  $E'$ , this means that given the typing declarations in  $\bar{\Gamma}, v : Y$ , the term  $E'$  must be assigned the type  $Z$ , and if  $v$  does not occur free in  $E'$ , this

means that given the typing declarations in  $\bar{\Gamma}$ , the term  $E'$  must be assigned the type  $Z$ . So by inductive hypothesis we have

$$\bar{\Gamma}, v : Y \vdash E' : Z \quad \text{or} \quad \bar{\Gamma} \vdash E' : Z$$

The same argument as in Theorem 1.43 then shows that in either case we have

$$\bar{\Gamma} \vdash \lambda v^Y(E') : Y \rightarrow Z$$

as desired, and **(i)**  $\rightarrow$  **(ii)** holds in this case. This completes the proof that **(i)**  $\rightarrow$  **(ii)** holds in general.

The proof that **(ii)**  $\rightarrow$  **(i)** is by induction on the construction of proofs. If  $\Gamma \vdash E : X$  has the form of an initial sequent  $\Delta, v : X \vdash v : X$ , then the result is trivial.

Suppose then that  $\Gamma \vdash E_1(E_2) : X$  is the result of an (App) inference from

$$\Gamma \vdash E_1 : Y \rightarrow X \quad \text{and} \quad \Gamma \vdash E_2 : Y$$

where  $v_1, \dots, v_n$  are the free variables of  $E_1(E_2)$ . Thus,  $\Gamma$  will include the typing declarations  $v_1 : \tau_1, \dots, v_n : \tau_n$  for some  $\tau_1, \dots, \tau_n$ . As before, let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, \dots, v_n : \tau_n$ , let  $\bar{\Gamma}_1$  be the subset of variable declarations in  $\bar{\Gamma}$  for free variables of  $E_1$ , and let  $\bar{\Gamma}_2$  be the subset of variable declarations in  $\bar{\Gamma}$  for free variables of  $E_2$ .

By inductive hypothesis, the variable declarations in  $\bar{\Gamma}_1$  entail that  $E_1$  is assigned the type  $Y \rightarrow X$ , and the variable declarations in  $\bar{\Gamma}_2$  entail that  $E_2$  is assigned the type  $Y \rightarrow X$ . Because  $\bar{\Gamma}_1$  and  $\bar{\Gamma}_2$  are contained in  $\bar{\Gamma}$ , it follows that the variable declarations in  $\bar{\Gamma}$  entail that  $E_1(E_2)$  is assigned the type  $X$ , as desired.

Suppose finally that  $\Gamma \vdash \lambda v^X(E') : X \rightarrow Y$  is the result of an (Abs<sup>bt</sup>) inference from

$$\Gamma, v : X \vdash E' : Y$$

where if  $v_1, \dots, v_n$  are the free variables of  $\lambda v^X(E')$ , then  $\Gamma$  includes the typing declarations  $v_1 : \tau_1, \dots, v_n : \tau_n$ . Let  $\bar{\Gamma}$  be the context  $v_1 : \tau_1, \dots, v_n : \tau_n$ .

If  $v$  occurs free in  $E'$ , then by inductive hypothesis, the variable declarations in  $\bar{\Gamma}, v : X$  require  $E'$  to be assigned the type  $Y$ . This means that the variable declarations in  $\bar{\Gamma}$  require  $\lambda v^X(E')$  to be assigned the type  $X \rightarrow Y$ , as desired.

If  $v$  does not occur free in  $E'$ , then by inductive hypothesis, the variable declarations in  $\bar{\Gamma}$  require  $E'$  to be assigned the type  $Y$ . This also means that the variable declarations in  $\bar{\Gamma}$  require  $\lambda v^X(E')$  to be assigned the type  $X \rightarrow Y$ , as desired. This completes the proof that **(ii)**  $\rightarrow$  **(i)** holds in general.

## Chapter 2

# Intuitionistic Propositional Logic and the Curry-Howard Correspondence.

### 2.1 Preamble

In this chapter, our main goal is to understand the Curry-Howard correspondence in its simplest forms. Broadly speaking, the Curry-Howard correspondence connects the notions of computation and proof - two notions which might otherwise seem quite distinct. There are many different forms of the Curry-Howard correspondence; in this chapter, we examine a simple version that connects the notion of certain sorts of computation with the notion of provability in *intuitionistic propositional logic*. In later chapters we examine more sophisticated versions of the Curry-Howard correspondence that connect richer notions of computation with richer notions of provability. We begin with a brief discussion of the main ideas behind intuitionistic logic.

### 2.2 Intuitionistic Logic: Motivation

One of the guiding ideas behind intuitionistic logic is that our reasoning ought to be *constructive*. So for example, according to constructive reasoning the only situation in which we are entitled to assert a disjunction  $A \vee B$  would be a situation in which either we are entitled to assert  $A$  or we are entitled to assert  $B$ . To put this in terms of knowledge, we could say that we only know a disjunctive claim  $A \vee B$  if we know  $A$  or know  $B$ . Intuitionistic logic makes similar claims about existentially quantified sentences  $(\exists x)\phi(x)$ . The only situation in which we are entitled to assert such a sentence is a situation in which we can name (or construct) an object  $o$  for which  $\phi(o)$ . To put it in terms of knowledge, we only know a claim of the form  $(\exists x)\phi(x)$  if we know that  $\phi(o)$  for some particular  $o$ .

Intuitionistic logic places different demands on us than classical logic. This point is illustrated nicely by the following classic example:

**Theorem 2.1**

There are irrational numbers  $x, y$  such that  $x^y$  is rational.

**Proof**

Consider the number  $(\sqrt{2}^{\sqrt{2}})^{\sqrt{2}}$ . Using the fact that  $(a^b)^c = a^{(b \cdot c)}$ , this equals 2. So if  $\sqrt{2}^{\sqrt{2}}$  is irrational, then  $x = \sqrt{2}^{\sqrt{2}}, y = \sqrt{2}$  gives an example of irrational numbers  $x$  and  $y$  such that  $x^y$  is rational. On the other hand, if  $\sqrt{2}^{\sqrt{2}}$  is rational, then  $x = \sqrt{2}, y = \sqrt{2}$  gives an example of irrational numbers  $x$  and  $y$  such that  $x^y$  is rational. So either way, there are irrational numbers  $x, y$  such that  $x^y$  is rational.

From the point of view of classical logic, this proof is completely persuasive. From the point of view of intuitionistic logic, it is not a valid proof at all. The theorem being proved is of course an existential claim; it states that *there are* irrational numbers  $x, y$  such that  $x^y$  is rational. But at the end of the proof we are still not in a position to name specific irrational numbers  $x$  and  $y$  for which  $x^y$  is rational. The intuitionistic logician is bothered by this - what could it mean to have proved an existential claim if a confirming instance has not been unambiguously specified? If one thinks that to prove a disjunction one must prove one of the disjuncts, or if one thinks that to prove an existential claim one must construct a specific confirming instance, then the above proof is completely unpersuasive.

Intuitionistic logic is an attempt to restrict our techniques of reasoning in such a way that whenever we can prove a disjunction  $A \vee B$ , then we know which of  $A$  or  $B$  obtains, and whenever we can prove an existential claim  $(\exists x)\phi(x)$ , then we can name an object  $o$  for which we have  $\phi(o)$ . It is perhaps surprising that one can quite formally and rigorously define a system of reasoning that conforms to these principles, and can in fact do so merely by making some small modifications to the usual rules of classical logic. Exploring how this is so will be one of the main themes of this chapter.

Although intuitionistic logic has a propositional form (in which there are no quantifiers) and a predicate form (which includes the usual quantifiers), in this chapter we will restrict our attention to propositional intuitionistic logic, turning to predicate intuitionistic logic in a later chapter.

First, we have the following definition

**Definition 2.2:**

In propositional intuitionistic logic, all formulae are constructed from atomic formulae and the the 0-ary connective  $\perp$  (sometimes called ‘the false’ or ‘falsum’) using the binary connectives  $\rightarrow, \vee$  and  $\&$ .

In what follows,  $\neg\phi$  will be an abbreviation for  $\phi \rightarrow \perp$ ,  $\phi \leftrightarrow \psi$  will be an abbreviation for  $(\phi \rightarrow \psi) \& (\psi \rightarrow \phi)$ , and  $\top$  (sometimes called ‘the true’ or ‘verum’) will be an abbreviation for  $\perp \rightarrow \perp$ .

So for example

$(\neg A) \vee B$	is an abbreviation of	$(A \rightarrow \perp) \vee B$
$(\neg A) \leftrightarrow B$	is an abbreviation of	$((A \rightarrow \perp) \rightarrow B) \& (B \rightarrow (A \rightarrow \perp))$
$(\neg A) \vee \neg \top$	is an abbreviation of	$(A \rightarrow \perp) \vee ((\perp \rightarrow \perp) \rightarrow \perp)$

Next we turn to a discussion of the connectives themselves. The logical constant  $\perp$  and the connectives  $\rightarrow, \vee$  and  $\&$  are the basic connectives of propositional intuitionistic logic. How are we to analyse their meaning? In the usual presentations of classical propositional logic, much emphasis is placed on *truth tables*. In particular, the connectives are often regarded as defined by the way they transform truth values. For example, consider the truth table for  $\&$ :

A	B	A&B
$\top$	$\top$	$\top$
$\top$	$\perp$	$\perp$
$\perp$	$\top$	$\perp$
$\perp$	$\perp$	$\perp$

This table presents  $\&$  as if it were a function taking a pair of truth values as inputs (given in the first two columns), and outputting a truth value (given in the third column). From the point of classical logic, one can even think of this as something like the *definition* of the connective  $\&$  – it is the binary connective that transforms truth values in the way shown in this table.

Intuitionistic logic takes a different point of view in analysing the meaning of the connectives. Rather than analysing the connectives in terms of truth values, intuitionistic logic analyzes them in terms of *proofs*, or even more broadly in terms of *constructions*. The basic ideas intuitionistic logic begins with are then as follows:

**Definition 2.3: The BHK Conditions**

**BHK1** : To prove  $A \& B$  is to prove  $A$  and to prove  $B$ .

**BHK2** : To prove  $A \vee B$  is to prove  $A$  or to prove  $B$ .

**BHK3** : To prove  $A \rightarrow B$  is to construct a way of transforming a proof of  $A$  into a proof of  $B$ .  
**BHK4** : There is no proof of  $\perp$ .

These conditions are described here somewhat informally, but we will later consider much more precise versions of them. They are sometimes described as the **BHK** conditions on the connectives in honor of Brouwer, Heyting, and Kolmogorov who first developed many of the main ideas of intuitionistic logic.

Consider **BHK1**. This does not just say that one way of proving  $A \& B$  is to prove  $A$  and prove  $B$ ; it rather makes the stronger claim that *what it is* to prove  $A \& B$  is to prove  $A$  and to prove  $B$ . Moreover, this is taken as our basic analysis of  $\&$ . That is to say, the *meaning* of  $\&$  is constituted by the fact that it is the connective having the property that proving  $A \& B$  just amounts to proving  $A$  and proving  $B$ .

Similarly, according to **BHK2**, the meaning of  $\vee$  is constituted by the fact that it is the connective having the property that proving  $A \vee B$  amounts to either a proof of  $A$  or a proof of  $B$ .

The idea behind **BHK3** is that to prove  $A \rightarrow B$  is to have some method for transforming a proof of  $A$  into a proof of  $B$ ; this is sometimes expressed by the idea that a proof of  $A \rightarrow B$  is a *function* that transforms a proof of  $A$  into a proof of  $B$ . This is usually more specifically taken to mean that there is some kind of *mechanical procedure* for transforming a proof of  $A$  into a proof of  $B$ . (At this point, we can be somewhat open as to what counts as a ‘mechanical procedure’.) The connective  $\rightarrow$  can even be regarded as being defined by the fact that it is the connective having the property that proving  $A \rightarrow B$  amounts to specifying a function that takes a proof of  $A$  into a proof of  $B$ .

The claim **BHK4** is perhaps more modest in character than **BHK1-3**. Presumably, we do *not* want to say that  $\perp$  is characterized by the fact that there is no proof of  $\perp$ , because there are sentences other than  $\perp$  which are not provable. While it might be argued that **BHK4** captures *part* of the meaning of  $\perp$ , it certainly does not capture all of it. Nevertheless, the claim made in **BHK4** is a true one, and it will suffice in what follows to regard **BHK4** as a basic truth about  $\perp$ .

The way of thinking of the connectives captured in **BHK1-4** is known as the **BHK** (Brouwer-Heyting-Kolmogorov) interpretation of intuitionistic logic. Note that according to this interpretation, because we regard  $\neg A$  as an abbreviation of  $A \rightarrow \perp$ , proving  $\neg A$  amounts to giving a procedure that takes a proof of  $A$  into a proof of  $\perp$  (i.e., a proof of a contradiction.)

In **BHK1-4**, the main non-classical idea lies in **BHK2** – that is, in the claim that proving  $A \vee B$  amounts to proving  $A$  or proving  $B$ . The classical logician could in fact happily accept **BHK1**, **BHK3**, and **BHK4**, but would emphatically reject **BHK2**. This is because in classical logic, one can have a proof of  $A \vee B$  without having either a proof of  $A$  or a proof of  $B$ . We have seen this already in our discussion of Theorem 2.1. For an even simpler example, consider the claim

**Even**: There are an even number of electrons in the universe.

Using classical logic, we can easily prove

**Even**  $\vee$   $\neg$ **Even**,

as this is a tautology. This however does not require that we have a proof of **Even** or a proof of  $\neg$ **Even**. Only scientific investigation can determine whether **Even** is true or false; logic on its own gives us neither a proof of **Even** nor a proof of  $\neg$ **Even**. The claim that proving  $A \vee B$  amounts to proving  $A$  or proving  $B$  is therefore a departure from classical logic. In effect, it imposes much stricter requirements on proving a disjunction  $A \vee B$  than are imposed by classical logic. According to the **BHK** conditions, we do *not* have a proof of **Even**  $\vee$   $\neg$ **Even, and perhaps never will.**

We have not yet spelt out any sort of formal system that captures the ideas of **BHK1-4**. However, even without doing so, we can see somewhat informally that various sentences should be theorems of intuitionistic logic. For example, consider the sentence  $(p \& q) \rightarrow p$ . This is a classical tautology. Should we expect this to be intuitionistically provable? According to **BHK3**, a proof of  $(p \& q) \rightarrow p$  would be a procedure that transforms a proof of  $p \& q$  into a proof of  $p$ . Given that according to **BHK1** a proof of a  $p \& q$  is just a proof of  $p$  and a proof of  $q$ , an intuitionistic proof of  $(p \& q) \rightarrow p$  would just be a function that takes a proof of  $p$  and a proof of  $q$  and returns a proof of  $p$ . Such a function, of course, exists - given a proof of  $p$  and a proof of  $q$ , just forget about the proof of  $q$ . One then has a proof of  $p$ . Thus,  $(p \& q) \rightarrow p$  should be intuitionistically provable.

Consider the sentence  $\perp \rightarrow p$ . This too is a classical tautology - anything follows from  $\perp$ . Should we expect this to be intuitionistically provable? For  $\perp \rightarrow p$  to be provable would be to have a way of transforming a proof of  $\perp$  into a proof of  $p$ . Because there are no proofs of  $\perp$ , we think of ourselves as 'vacuously' being able to transform a proof of  $\perp$  into a proof of  $p$  - that is to say, any transformational procedure may be said to take a proof of  $\perp$  into a proof of  $p$ , because there are no proofs of  $\perp$ . So this should be intuitionistically provable.

Consider now the classically tautology  $p \rightarrow (q \rightarrow p)$ . Should we expect this to be intuitionistically provable? The question we must ask is whether we have a way of transforming a proof of  $p$  into way of transforming a proof of  $q$  into a proof of  $p$ . Indeed, we do. Given a proof  $\Pi$  of  $p$ , we have a way of transforming a proof  $\Sigma$  of  $q$  into a proof of  $p$  as follows - erase  $\Sigma$ , and write down  $\Pi$ . This transforms any proof of  $q$  into a proof of  $p$ . Thus, we should expect  $p \rightarrow (q \rightarrow p)$  to be intuitionistically provable.

Likewise, you should be able to convince yourself that  $p \rightarrow \neg\neg p$  ought to be intuitionistically provable. (This is included in the exercises.)

There are of course classical tautologies that we would *not* expect to be intuitionistically provable. We have already considered the sentence  $p \vee \neg p$ . For this to be provable is for either  $p$  or  $\neg p$  to be provable. Of course, it is not the case that for every sentence  $p$ , we have that either  $p$  or  $\neg p$  is provable using pure logic. So we should expect that  $p \vee \neg p$  will *not* be intuitionistically provable in general. In this sense, the law of the excluded middle fails for intuitionistic logic, as a result of **BHK2**. You should also try to convince yourself that there are no good reasons to expect  $\neg\neg p \rightarrow p$  be intuitionistically provable. Perhaps a little more surprisingly, you should be able to convince yourself that there are no good reasons to expect *Peirce's Law*  $((p \rightarrow q) \rightarrow p) \rightarrow p$  to be intuitionistically provable, even though this is a purely implicational classical tautology and does not involve negation or disjunction in any way.

Of course, none of the arguments here suggesting that something should or should not be intuitionistically provable are completely rigorous. They are simply argument sketches. Armed

with further material, we will later be able to replace such argument sketches with completely rigorous arguments. Arguments sketches of the sort we have considered here are nevertheless suggestive and have played an important historical role in the development of intuitionistic logic.

### Exercises for Section 2.2

1. Using the **BHK** interpretation, argue informally that the following should be intuitionistically provable. It suffices to give informal arguments of the sort given in this section of the text.

- (a)  $p \rightarrow (p \vee q)$
- (b)  $p \rightarrow \neg\neg p$
- (c)  $(p \rightarrow q) \rightarrow (\neg q \rightarrow \neg p)$
- (d)  $\neg\neg\neg p \rightarrow \neg p$

2. Using the **BHK** interpretation, argue informally that the following should not be intuitionistically provable. It suffices to describe how naive ways of trying to prove them using the **BHK** rules fail.

- (a)  $\neg\neg p \rightarrow p$
- (b)  $(\neg q \rightarrow \neg p) \rightarrow (p \rightarrow q)$
- (c)  $((p \rightarrow q) \rightarrow p) \rightarrow p$

## 2.3 Natural Deduction for Propositional Intuitionistic Logic

We now set up a formal proof system that captures the informal notion of intuitionistic provability developed in the last section. It will not be clear until somewhat later in this chapter why or how the proof system presented in this section captures the main ideas behind the **BHK** interpretation. In particular, only once we understand the Curry-Howard correspondence and proof normalization will we be able to clearly see the sense in which this proof system captures the **BHK** rules. Until then, the reader should just focus on the details of the proof system itself.

In our proof system, each step in the proof will consist of a sequent of the form

$$s_1, s_2, \dots, s_n \vdash t$$

where  $s_1, s_2, \dots, s_n, t$  are formulae. For brevity, we will typically write this as

$$\Delta \vdash t$$

where  $\Delta$  is the unordered list of sentences  $s_1, s_2, \dots, s_n$ . In  $\Delta$ , we will allow elements to be repeated. Intuitively, the sequent  $\Delta \vdash t$  says that  $t$  is derivable from the sentences in  $\Delta$ . We will write

$$\vdash s$$



as an abbreviation of

$$\emptyset \vdash s,$$

where  $\emptyset$  is the empty list. To say  $\vdash s$  is then to say that  $s$  is provable outright.

A proof will consist of a ‘tree’ of such sequents. The leaves of the trees will be located at the top, and will correspond to our logical axioms. The statement being proven will lie at the base of the tree. Often we are interested in proving a statement of the form  $\vdash s$ , but we will allow more general assertions of the form  $\Delta \vdash s$  to be the conclusion of a proof and lie at the base of the tree. We will get from one level of a tree to the level below by performing a rule of inference. Examples will be provided shortly.

We specify the logical axioms and rules of inference of our system, which we call **NJ**.

**Definition 2.4: The Intuitionistic Propositional Calculus NJ.**

For any formulae  $\phi, \psi, \theta$  and any finite unordered list of formulae  $\Delta$ , we have

**Logical Axioms:**

$$\Gamma, \phi \vdash \phi$$

**Rules of Inference:**

$$\frac{\Delta, \phi \vdash \psi}{\Delta \vdash \phi \rightarrow \psi} (\rightarrow I) \quad \frac{\Delta \vdash \phi \rightarrow \psi \quad \Delta \vdash \phi}{\Delta \vdash \psi} (\rightarrow E)$$

(major premise:  $\Delta \vdash \phi \rightarrow \psi$ )

$$\frac{\Delta \vdash \phi \quad \Delta \vdash \psi}{\Delta \vdash \phi \& \psi} (\& I) \quad \frac{\Delta \vdash \phi \& \psi}{\Delta \vdash \phi} (\& E) \quad \frac{\Delta \vdash \phi \& \psi}{\Delta \vdash \psi} (\& E)$$

$$\frac{\Delta \vdash \phi}{\Delta \vdash \phi \vee \psi} (\vee I) \quad \frac{\Delta \vdash \psi}{\Delta \vdash \phi \vee \psi} (\vee I)$$

$$\frac{\Delta \vdash \phi \vee \psi \quad \Delta, \phi \vdash \theta \quad \Delta, \psi \vdash \theta}{\Delta \vdash \theta} (\vee E)$$

(major premise:  $\Delta \vdash \phi \vee \psi$ )

$$\frac{\Delta \vdash \perp}{\Delta \vdash \phi} (\perp E)$$

In this system when something like  $\Delta, \phi$  appears on the left hand side of a  $\vdash$  symbol, we understand  $\Delta, \phi$  as the unordered list of sentences obtained by adding  $\phi$  to the unordered list of sentences  $\Delta$ .

In the case of the  $(\rightarrow E)$  rule, we say that the premise  $\Delta \vdash \phi \rightarrow \psi$  is the *major* premise, and the remaining premise is *minor*, and in the case of the  $(\vee E)$  rule, we say that the premise  $\Delta \vdash \phi \vee \psi$  is the major premise, and the remaining premises minor. For the remaining rules, we will not need to classify premises as major or minor.

Note that the rules and logical axioms of **NJ** are all classically acceptable. Unlike classical logic, however, **NJ** does *not* have the following rule:

$$\frac{\Delta, \neg\phi \vdash \perp}{\Delta \vdash \phi} \quad (\times)$$

This rule is of course classically valid, but we do not include it in our system, and it is not derivable in any way from the rules of **NJ**.

It is not obvious that this moderate restriction of classical logic captures the ideas behind the **BHK** interpretation, but we will see quite rigorously by the end of the chapter that it does. Traditional treatments of intuitionistic logic develop a semantics for intuitionistic logic (typically using the machinery of ‘Kripke frames’), and then prove a soundness and completeness theorem relative to this semantics, in analogy with the familiar soundness and completeness theorems of classical logic. (See [7] and chapter 2 of [10] for examples of this.) We will not take this route. Rather than thinking of intuitionistic proofs as arguments that preserve truth in Kripke frames, we will end up thinking of intuitionistic proofs as computer programs or algorithms (given by typed  $\lambda$ -terms) that obey the **BHK** rules.

As simple examples of proofs in intuitionistic logic, here are intuitionistic proofs of  $p \rightarrow (q \rightarrow p)$  and  $p \rightarrow \neg\neg p$ , recalling that this latter sentence is an abbreviation for  $p \rightarrow ((p \rightarrow \perp) \rightarrow \perp)$ :

$$\frac{\frac{p, q \vdash p}{p \vdash q \rightarrow p} (\rightarrow \text{I})}{\vdash p \rightarrow (q \rightarrow p)} (\rightarrow \text{I})$$

$$\frac{\frac{\frac{p, p \rightarrow \perp \vdash p}{p, p \rightarrow \perp \vdash \perp} (\rightarrow \text{E})}{p \vdash (p \rightarrow \perp) \rightarrow \perp} (\rightarrow \text{I})}{\vdash p \rightarrow ((p \rightarrow \perp) \rightarrow \perp)} (\rightarrow \text{I})$$

The classical principle that everything follows from a contradiction also holds in intuitionistic logic:

$$\frac{\frac{\frac{p \& \neg p \vdash p \& \neg p}{p \& \neg p \vdash p} (\& \text{E}) \quad \frac{p \& \neg p \vdash p \& \neg p}{p \& \neg p \vdash \neg p} (\& \text{E})}{\frac{p \& \neg p \vdash \perp}{p \& \neg p \vdash q} (\perp \text{E})}{\vdash (p \& \neg p) \rightarrow q} (\rightarrow \text{I})$$

For a more complicated argument, note that using the proofs

$$\frac{\frac{(\neg p) \vee (\neg q), p \& q, \neg p \vdash \neg p}{(\neg p) \vee (\neg q), p \& q, \neg p \vdash \perp} (\& \text{I})}{(\neg p) \vee (\neg q), p \& q, \neg p \vdash \perp} (\rightarrow \text{E})$$

and

$$\frac{(\neg p) \vee (\neg q), p \& q, \neg q \vdash \neg q \quad \frac{(\neg p) \vee (\neg q), p \& q, \neg q \vdash p \& q}{(\neg p) \vee (\neg q), p \& q, \neg q \vdash q} (\& \text{I})}{(\neg p) \vee (\neg q), p \& q, \neg q \vdash \perp} (\rightarrow \text{E})$$

we can prove  $((\neg p) \vee (\neg q)) \rightarrow \neg(p \& q)$  :

$$\frac{(\neg p) \vee (\neg q), p \& q \vdash (\neg p) \vee (\neg q) \quad \frac{\vdots \quad \frac{(\neg p) \vee (\neg q), p \& q, \neg p \vdash \perp}{(\neg p) \vee (\neg q), p \& q \vdash \perp} (\rightarrow \text{I}) \quad \frac{\vdots \quad \frac{(\neg p) \vee (\neg q), p \& q, \neg q \vdash \perp}{(\neg p) \vee (\neg q), p \& q \vdash \perp} (\rightarrow \text{I})}{(\neg p) \vee (\neg q) \vdash \neg(p \& q)} (\rightarrow \text{I})}{\vdash ((\neg p) \vee (\neg q)) \rightarrow \neg(p \& q)} (\rightarrow \text{I}) (\vee \text{E})$$

We have called the proof system just presented **NJ**. The system **NJ**( $\rightarrow$ ) denotes the ‘implicational’ fragment of this system - that is, the system with the same logical axioms, but with only the rules ( $\rightarrow$  I) and ( $\rightarrow$  E).

**Definition 2.5: The System NJ( $\rightarrow$ ).**

**Logical Axioms:**

$$\Gamma, \phi \vdash \phi$$

**Rules of Inference:**

$$\frac{\Delta, \phi \vdash \psi}{\Delta \vdash \phi \rightarrow \psi} (\rightarrow \text{I}) \quad \frac{\Delta \vdash \phi \rightarrow \psi \quad \Delta \vdash \phi}{\Delta \vdash \psi} (\rightarrow \text{E})$$

(major premise:  $\Delta \vdash \phi \rightarrow \psi$ )

Of the many familiar theorems from classical logic, **NJ** proves some but not others. Most notably, it does not prove the law of the excluded middle  $p \vee \neg p$  in general. (Later in this chapter, we will rigorously *prove* that **NJ** does not prove the law of the excluded middle!) And although we have seen that **NJ** can prove  $p \rightarrow \neg\neg p$ , it cannot prove  $\neg\neg p \rightarrow p$  in general. (Interestingly, it can prove  $\neg\neg\neg p \rightarrow \neg p$ .) Not all the familiar deMorgan’s laws hold either. So although the following

$$\begin{aligned} (\neg p \& \neg q) &\rightarrow \neg(p \vee q) \\ \neg(p \vee q) &\rightarrow (\neg p \& \neg q) \\ (\neg p \vee \neg q) &\rightarrow \neg(p \& q). \end{aligned}$$

are all provable in **NJ**,

$$\neg(p \& q) \rightarrow (\neg p \vee \neg q)$$

is not provable in **NJ**. However, all the familiar distributivity rules relating conjunction and disjunction

$$(p \vee (q \& r)) \rightarrow ((p \vee q) \& (p \vee r))$$

$$\begin{aligned} &((p \vee q) \& (p \vee r)) \rightarrow (p \vee (q \& r)) \\ &((p \& q) \vee (p \& r)) \rightarrow (p \& (q \vee r)) \\ &(p \& (q \vee r)) \rightarrow ((p \& q) \vee (p \& r)) \end{aligned}$$

are provable in **NJ**. In the exercises you will verify many of these facts about what is provable in **NJ**. Later in this chapter we will be in a position to show rigorously that certain sentences are *not* provable in **NJ**.

### Exercises for Section 2.3

1. Prove the following sentences in **NJ**:

- (a)  $p \rightarrow (p \vee q)$
- (b)  $\neg(p \vee q) \rightarrow \neg p$
- (c)  $\neg(p \vee q) \rightarrow (\neg p) \& (\neg q)$
- (d)  $(\neg p) \& (\neg q) \rightarrow \neg(p \vee q)$
- (e)  $(\neg p) \vee (\neg q) \rightarrow \neg(p \& q)$
- (f)  $p \rightarrow \neg\neg p$
- (g)  $(p \vee (q \& r)) \rightarrow ((p \vee q) \& (p \vee r))$
- (h)  $((p \vee q) \& (p \vee r)) \rightarrow (p \vee (q \& r))$
- (i)  $((p \& q) \vee (p \& r)) \rightarrow (p \& (q \vee r))$
- (j)  $(p \& (q \vee r)) \rightarrow ((p \& q) \vee (p \& r))$
- (k)  $\neg\neg(p \vee \neg p)$
- (l)  $(p \rightarrow q) \rightarrow (\neg q \rightarrow \neg p)$
- (m)  $\neg\neg\neg p \rightarrow \neg p$

2. Show that a formula  $X$  is provable in classical propositional logic iff  $\neg\neg X$  is provable in **NJ** as follows. (This result is sometimes called *Glivenko's Theorem*.) For our purposes, we take classical propositional logic to be **NJ** with additional logical axioms of the form

$$\Delta \vdash X \vee \neg X$$

permitted for any sentence  $X$  and unordered list of sentences  $\Delta$ .

(a) Prove that if  $\vdash \neg\neg X$  is provable in **NJ**, then  $\vdash X$  is provable in classical propositional logic. (Hint: prove and use the fact that  $\neg\neg X \vdash X$  is derivable in classical propositional logic, using  $(\vee E)$ .)

(b) Show that if  $\vdash X$  is provable in classical propositional logic, then for some formulae  $p_1, \dots, p_n$

$$p_1 \vee \neg p_1, p_2 \vee \neg p_2, \dots, p_n \vee \neg p_n \vdash X$$

is provable in **NJ**. Infer from this that

$$p_1 \vee \neg p_1, p_2 \vee \neg p_2, \dots, p_n \vee \neg p_n, \neg X \vdash \perp$$

is provable in **NJ**.

(c) From this, show that if  $\vdash X$  is provable in classical propositional logic, then  $\vdash \neg\neg X$  is provable in **NJ**. (Hint: use the fact that  $\vdash \neg\neg(p \vee \neg p)$  is provable in **NJ**.)

## 2.4 Some Abbreviations

It will help to introduce some conventions to simplify the presentation of proofs. The following rule, sometimes called weakening, allows us to introduce an additional premise on the left of a sequent

$$\frac{\Delta \vdash B}{\Delta, A \vdash B} \text{ (W)}$$

In fact, anything we can prove with this rule we can prove without it. For example, consider the proof

$$\frac{\frac{B, B \rightarrow A \vdash B \quad B, B \rightarrow A \vdash B \rightarrow A}{B, B \rightarrow A \vdash A} (\rightarrow E)}{B, B \rightarrow A, C \vdash A} \text{ (W)} \quad (31)$$

In the final step of this proof, we introduce  $C$  on the left by a weakening. However, we could in fact have just included  $C$  in the left hand side of each sequent all along, thereby eliminating the need for a weakening, as follows:

$$\frac{B, B \rightarrow A, C \vdash B \quad B, B \rightarrow A, C \vdash B \rightarrow A}{B, B \rightarrow A, C \vdash A} (\rightarrow E) \quad (32)$$

Although weakening adds an extra step to the proof, it allows us to introduce a formula on the left hand side of a sequent only when it is needed, rather than being forced to carry it around from the beginning of a proof. Weakening can thus make our proofs more compact, and so we allow ourselves to use this rule. However, we do not regard weakening as one of the basic rules of **NJ**. Instead, we think of (31) as an ‘abbreviation’ of (32). So (31) is not technically a proof of **NJ**, but rather an abbreviation of proof (32), which is a proof of **NJ**.

It will be useful to introduce several other ‘abbreviations’ of this sort. Consider the rule ( $\rightarrow E$ )

$$\frac{\Delta \vdash A \rightarrow B \quad \Delta \vdash A}{\Delta \vdash B} (\rightarrow E)$$

In this rule the context  $\Delta$  on the left hand side of each sequent is required to be the same. An alternative would be to not require this, but instead require that the context in the conclusion be the concatenation of the contexts in the premises. We call this rule ( $\rightarrow E'$ )

$$\frac{\Delta \vdash A \rightarrow B \quad \Delta' \vdash A}{\Delta, \Delta' \vdash B} (\rightarrow E')$$

where by  $\Delta, \Delta'$  we simply mean the unordered list consisting of  $\Delta$  conjoined with  $\Delta'$ . If  $\Delta$  and  $\Delta'$  have an element  $X$  in common, then  $X$  may of course appear multiple times in  $\Delta, \Delta'$ . We allow

the option of eliminating as many such ‘duplicate’ formulae as we want on the left hand side of the sequent when applying this rule. So for example, each of the following are valid instances of ( $\rightarrow E'$ )

$$\frac{A, A \rightarrow B, C_1 \vdash A \rightarrow B \quad A, A \rightarrow B, C_2 \vdash A}{A, A, A \rightarrow B, A \rightarrow B, C_1, C_2 \vdash B} (\rightarrow E')$$

$$\frac{A, A \rightarrow B, C_1 \vdash A \rightarrow B \quad A, A \rightarrow B, C_2 \vdash A}{A, A \rightarrow B, A \rightarrow B, C_1, C_2 \vdash B} (\rightarrow E') \quad (33)$$

$$\frac{A, A \rightarrow B, C_1 \vdash A \rightarrow B \quad A, A \rightarrow B, C_2 \vdash A}{A, A \rightarrow B, C_1, C_2 \vdash B} (\rightarrow E')$$

Like the rule of weakening, we do not need to think of ( $\rightarrow E'$ ) as a new basic rule of **NJ**. In fact, anything we can prove with ( $\rightarrow E'$ ) can be proven just with ( $\rightarrow E$ ) by adding needed formulae to the left hand side of sequents above the ( $\rightarrow E$ ), in much the same way that we did for weakening. For example, the middle proof (33) can be regarded as an abbreviation for

$$\frac{A, A \rightarrow B, A \rightarrow B, C_1, C_2 \vdash A \rightarrow B \quad A, A \rightarrow B, A \rightarrow B, C_1, C_2 \vdash A}{A, A \rightarrow B, A \rightarrow B, C_1, C_2 \vdash B} (\rightarrow E) \quad (34)$$

In this way, we regard proofs using ( $\rightarrow E'$ ) as simply abbreviations of proofs using the conventional ( $\rightarrow E$ ) rule. Comparing (33) and (34) we see that the ( $\rightarrow E'$ ) rule can simplify our proofs, and for this reason we use it.

Using these new abbreviations, we can compress (31) even further as follows

$$\frac{\frac{B \vdash B \quad B \rightarrow A \vdash B \rightarrow A}{B, B \rightarrow A \vdash A} (\rightarrow E')}{B, B \rightarrow A, C \vdash A} (W) \quad (35)$$

We allow similar ‘abbreviations’ for all our rules which take more than one sequent as premise – so not just ( $\rightarrow E$ ), but also ( $\&I'$ ) and ( $\vee E'$ ).

The material of this section is captured by the following theorem

**Theorem 2.6: The Admissibility of (W), ( $\rightarrow E'$ ), ( $\&I'$ ) and ( $\vee E'$ ).**

In each of the rules

$$\frac{\Delta \vdash B}{\Delta, A \vdash B} (W) \quad \frac{\Delta \vdash A \rightarrow B \quad \Delta' \vdash A}{\Delta, \Delta' \vdash B} (\rightarrow E')$$

$$\frac{\Delta \vdash A \quad \Delta' \vdash B}{\Delta, \Delta' \vdash A \& B} (\&I) \quad \frac{\Delta \vdash A \vee B \quad \Delta', A \vdash C \quad \Delta'', B \vdash C}{\Delta, \Delta', \Delta'' \vdash C} (\vee E)$$

if the premise(s) are provable in **NJ**, then so is the conclusion. In the conclusions of the rules ( $\rightarrow E'$ ), ( $\&I'$ ) and ( $\vee E'$ ), we allow the optional elimination of ‘duplicate’ copies of formulae

created when the contexts  $\Delta, \Delta', \dots$  partially overlap.

This theorem tells us that even though (W), ( $\rightarrow E'$ ), ( $\&I'$ ) and ( $\vee E'$ ) are not rules of **NJ**, we may treat them as if they were. The proof is a straightforward induction on the construction of proofs, and a representative case is left to the exercises.

### Exercises for Section 2.4

1. To what proofs using only the rules of **NJ** do the following proofs correspond?

$$(i) \quad \frac{C \vdash C \quad \frac{B \vdash B}{A, B \vdash B} (W)}{A, B, C \vdash B \& C} (\&I')$$

$$(ii) \quad \frac{A \vee B \vdash A \vee B \quad B \vdash B \quad \frac{A \vdash A \quad \neg A \vdash \neg A}{A, \neg A \vdash \perp} (\rightarrow E')}{A \vee B, \neg A \vdash B} (\vee E')$$

2. By using the rules (W), ( $\rightarrow E'$ ), ( $\&I'$ ) and ( $\vee E'$ ) wherever possible, simplify the following proofs

$$(i) \quad \frac{\frac{A, A \rightarrow \perp \vdash A \quad A, A \rightarrow \perp \vdash A \rightarrow \perp}{A, A \rightarrow \perp \vdash \perp} (\rightarrow E)}{A \vdash (A \rightarrow \perp) \rightarrow \perp} (\rightarrow I)}{\vdash A \rightarrow ((A \rightarrow \perp) \rightarrow \perp)} (\rightarrow I)$$

(ii) the proof of  $((\neg p) \vee (\neg q)) \rightarrow \neg(p \& q)$  given in section 2.3.

3. Prove that if  $\Delta \vdash B$  is provable in **NJ**, then so is  $\Delta, A \vdash B$ . (Hint: use induction on the construction of proofs.)

4. Prove that if  $\Delta \vdash A$  and  $\Delta' \vdash B$  are provable in **NJ**, then so is  $\Delta, \Delta' \vdash A \& B$ . (Hint: use the result of the previous exercise.)

## 2.5 The Curry-Howard Correspondence: Version I

We now turn our attention back to typing derivations and the arguments of the formal system **TR<sub>0</sub>**. This is a system in which it can be rigorously shown that a given untyped  $\lambda$ -term can be assigned a given type. We have seen that it is a central tool in the Curry-style typed  $\lambda$ -calculus.

The arguments of **TR<sub>0</sub>** are strikingly similar to ordinary proofs in logic. For example, take the typing derivation in which it is shown that the term  $\lambda x(\lambda y(x(xy)))$  can be assigned the type  $(A \rightarrow A) \rightarrow (A \rightarrow A)$ :

$$\frac{x : A \rightarrow A, y : A \vdash x : A \rightarrow A \quad \frac{x : A \rightarrow A, y : A \vdash x : A \rightarrow A \quad x : A \rightarrow A, y : A \vdash y : A}{x : A \rightarrow A, y : A \vdash xy : A}}{x : A \rightarrow A, y : A \vdash x(xy) : A}}{x : A \rightarrow A \vdash \lambda y(x(xy)) : A \rightarrow A}}{\vdash \lambda x(\lambda y(x(xy))) : (A \rightarrow A) \rightarrow (A \rightarrow A)}$$

If in each judgment of the form  $z : B$  we simply discard the left hand side  $z$  (the  $\lambda$ -term) and keep the right hand side  $B$  (the type), the above proof becomes the following:

$$\frac{A \rightarrow A, A \vdash A \rightarrow A \quad \frac{A \rightarrow A, A \vdash A \rightarrow A \quad A \rightarrow A, A \vdash A}{A \rightarrow A, A \vdash A}}{A \rightarrow A \vdash A \rightarrow A}}{\vdash (A \rightarrow A) \rightarrow (A \rightarrow A)}$$

This is now an ordinary proof in propositional logic using only the rules ( $\rightarrow$  I) and ( $\rightarrow$  E), and is thus a proof in **NJ**( $\rightarrow$ ). This should not be surprising. Our typing derivations use only the rules:

$$\begin{array}{c} \Gamma, x : A \vdash x : A \text{ (Var)} \\ \frac{\Gamma \vdash x : A \rightarrow B \quad \Gamma \vdash y : A}{\Gamma \vdash xy : B} \text{ (App)} \\ \frac{\Gamma, x : A \vdash y : B}{\Gamma \vdash \lambda xy : A \rightarrow B} \text{ (Abs)} \end{array}$$

The process of discarding the  $\lambda$ -terms and keeping the types in these rules yields:

$$\frac{\frac{\frac{\bar{\Gamma}, A \vdash A}{\bar{\Gamma} \vdash A \rightarrow B} \quad \bar{\Gamma} \vdash A}{\bar{\Gamma} \vdash B}}{\bar{\Gamma}, A \vdash B}}{\bar{\Gamma} \vdash A \rightarrow B}$$

where  $\bar{\Gamma}$  is the list of types in  $\Gamma$ . This of course is precisely the set of logical axioms and rules of **NJ**( $\rightarrow$ ).

A striking consequence of this is as follows. Suppose a term  $t$  of the untyped  $\lambda$ -calculus is closed (i.e., has no free variables), and can be assigned the type  $\tau$ . According to the Soundness and Completeness Theorem of **TR**<sub>0</sub> (Theorem 1.43), there is then a proof of  $\vdash t : \tau$  in **TR**<sub>0</sub>. By deleting terms and keeping types in this proof, we can transform this typing derivation in **TR**<sub>0</sub> into a proof of  $\vdash \tau$  in **NJ**( $\rightarrow$ ). Thus, if a closed term of the untyped  $\lambda$ -calculus can be assigned the type  $\tau$ , then  $\tau$  must be a theorem of **NJ**( $\rightarrow$ ). So for example, it immediately follows that the type  $\neg\neg A \rightarrow A$  will *not* in general be the type of *any* closed  $\lambda$ -term, as  $\neg\neg A \rightarrow A$  is not a theorem of **NJ**( $\rightarrow$ ).





This sort of strategy can be applied quite generally. We have

**Theorem 2.8: The Curry-Howard Correspondence for  $\mathbf{TR}_0$  and  $\mathbf{NJ}(\rightarrow)$ , Part 2.**

If  $\mathbf{NJ}(\rightarrow)$  proves

$$A_1, \dots, A_n \vdash B$$

Then for any distinct variables  $x_1, \dots, x_n$ , there is some  $\lambda$ -term  $s$  such that  $\mathbf{TR}_0$  proves

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

**Proof**

The proof is a straightforward induction on the construction of proofs in  $\mathbf{NJ}(\rightarrow)$ . For the base case, consider a leaf

$$A_1, \dots, A_n, B \vdash B$$

then for any distinct variables  $x_1, \dots, x_n, y$  we have that

$$x_1 : A_1, \dots, x_n : A_n, y : B \vdash y : B$$

is a proof in  $\mathbf{TR}_0$ . This completes the base case.

Suppose next that a proof of  $A_1, \dots, A_n \vdash B$  has as its final step an application of  $(\rightarrow E)$ :

$$\frac{\begin{array}{c} \vdots \\ A_1, \dots, A_n \vdash C \rightarrow B \end{array} \quad \begin{array}{c} \vdots \\ A_1, \dots, A_n \vdash C \end{array}}{A_1, \dots, A_n \vdash B} (\rightarrow E)$$

and let  $x_1, \dots, x_n$  be distinct variables. By inductive hypothesis, for some  $\lambda$ -terms  $s_1$  and  $s_2$ , there are proofs of

$$x_1 : A_1, \dots, x_n : A_n \vdash s_1 : C \rightarrow B \quad \text{and} \quad x_1 : A_1, \dots, x_n : A_n \vdash s_2 : C$$

in  $\mathbf{TR}_0$ . Combining these, we then have the following proof in  $\mathbf{TR}_0$

$$\frac{\begin{array}{c} \vdots \\ x_1 : A_1, \dots, x_n : A_n \vdash s_1 : C \rightarrow B \end{array} \quad \begin{array}{c} \vdots \\ x_1 : A_1, \dots, x_n : A_n \vdash s_2 : C \end{array}}{x_1 : A_1, \dots, x_n : A_n \vdash s_1(s_2) : B}$$

as desired.

Finally, suppose that a proof of  $A_1, \dots, A_n \vdash B$  has as its final step an application of  $(\rightarrow E)$ . This means that  $B$  has the form  $C \rightarrow D$ :

$$\frac{\begin{array}{c} \vdots \\ A_1, \dots, A_n, C \vdash D \end{array}}{A_1, \dots, A_n \vdash C \rightarrow D} (\rightarrow E)$$

Let  $x_1, \dots, x_n, x_{n+1}$  be distinct variables. By inductive hypothesis, for some  $\lambda$ -term  $s$  there is a proof of

$$x_1 : A_1, \dots, x_n : A_n, x_{n+1} : C \vdash s : D$$

in  $\mathbf{TR}_0$ . We then have the following proof in  $\mathbf{TR}_0$

$$\frac{\begin{array}{c} \vdots \\ x_1 : A_1, \dots, x_n : A_n, x_{n+1} : C \vdash s : D \end{array}}{x_1 : A_1, \dots, x_n : A_n \vdash \lambda x_{n+1}(s) : C \rightarrow D} (\rightarrow E)$$

as desired.

The two parts of the Curry-Howard correspondence tell us that there is an intimate connection between proofs in  $\mathbf{NJ}(\rightarrow)$  and proofs in  $\mathbf{TR}_0$ . Every proof in  $\mathbf{TR}_0$  contains as its ‘skeleton’ a proof in  $\mathbf{NJ}(\rightarrow)$ , and every proof in  $\mathbf{NJ}(\rightarrow)$  is the ‘skeleton’ of some proof in  $\mathbf{TR}_0$ .

Let us return to questions about the typing of terms, and introduce some terminology. Given a type  $A$ , one might wonder where there is some  $\lambda$ -term  $s$  for which  $s$  can be assigned type  $A$ . This is equivalent to asking whether there is some  $s$  such that  $\vdash s : A$  is provable in  $\mathbf{TR}_0$ . This is sometimes called the *type inhabitation problem* for  $A$ . If  $\vdash s : A$ , we say that the term  $s$  ‘inhabits’ the type  $A$ , and that the type  $A$  is ‘inhabited’. One can then ask which types are inhabited. Theorems 2.7 and 2.8 tell us that a type  $A$  is inhabited iff  $A$  is a theorem of  $\mathbf{NJ}(\rightarrow)$ . So a question about types is perhaps surprisingly answered by invoking a concept from intuitionistic logic.

Not only can thinking about intuitionistic logic help us answer questions about the typing of  $\lambda$ -terms, but thinking about the typing of  $\lambda$ -terms can also help us answer questions about intuitionistic logic. For example, let us ask: is the sentence

$$(A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C)) \quad (36)$$

a theorem of  $\mathbf{NJ}(\rightarrow)$ ? By the Curry-Howard correspondence, this sentence is a theorem of  $\mathbf{NJ}(\rightarrow)$  iff there is a closed  $\lambda$ -term that can be given the type (36). Careful experimentation shows that there is such a term, namely

$$\lambda x(\lambda y(\lambda z(y(xz)))) \quad (37)$$

which may be typed as follows:

$$\lambda x^{A \rightarrow B}(\lambda y^{B \rightarrow C}(\lambda z^A(y(xz)))).$$

That the  $\lambda$ -term (37) may indeed be given the type (36) is shown explicitly by the following typing derivation

$$\begin{array}{c}
\frac{x:A \rightarrow B, y:B \rightarrow C, z:A \vdash y:B \rightarrow C \quad \frac{x:A \rightarrow B, y:B \rightarrow C, z:A \vdash x:A \rightarrow B \quad x:A \rightarrow B, y:B \rightarrow C, z:A \vdash z:A}{x:A \rightarrow B, y:B \rightarrow C, z:A \vdash xz:B}}{x:A \rightarrow B, y:B \rightarrow C, z:A \vdash y(xz):C} \\
\frac{x:A \rightarrow B, y:B \rightarrow C \vdash \lambda z(y(xz)):A \rightarrow C}{x:A \rightarrow B \vdash \lambda y(\lambda z(y(xz))):(B \rightarrow C) \rightarrow (A \rightarrow C)} \\
\frac{\quad}{\vdash \lambda x(\lambda y(\lambda z(y(xz)))):(A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C))}
\end{array}$$

Eliminating the terms and keeping the types in this typing derivation then yields a proof of (37) in  $\mathbf{NJ}(\rightarrow)$ :

$$\begin{array}{c}
\frac{A \rightarrow B, B \rightarrow C, A \vdash B \rightarrow C \quad \frac{A \rightarrow B, B \rightarrow C, A \vdash A \rightarrow B \quad A \rightarrow B, B \rightarrow C, A \vdash A}{A \rightarrow B, B \rightarrow C, A \vdash B}}{A \rightarrow B, B \rightarrow C, A \vdash C} \\
\frac{A \rightarrow B, B \rightarrow C \vdash A \rightarrow C}{A \rightarrow B \vdash (B \rightarrow C) \rightarrow (A \rightarrow C)} \\
\frac{\quad}{\vdash (A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C))}
\end{array}$$

In this way, the  $\lambda$ -term  $\lambda x(\lambda y(\lambda z(y(xz))))$  ‘encodes’ a proof of  $(A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C))$  in  $\mathbf{NJ}(\rightarrow)$ . Thinking of  $\lambda$ -terms as programs, proofs are then encoded by programs. Searching for a proof of a given sentence in  $\mathbf{NJ}(\rightarrow)$  is then equivalent to searching for a closed  $\lambda$ -term (i.e., a program) with a given type.

Finally, it is worth pointing out that the Curry-Howard Correspondence also holds between  $\mathbf{TR}_0^{\text{bt}}$  and  $\mathbf{NJ}(\rightarrow)$ . In particular, we have the following theorems

**Theorem 2.9: The Curry-Howard Correspondence for  $\mathbf{TR}_0^{\text{bt}}$  and  $\mathbf{NJ}(\rightarrow)$ , Part 1.**

If  $\mathbf{TR}_0^{\text{bt}}$  proves

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

then  $\mathbf{NJ}(\rightarrow)$  proves

$$A_1, \dots, A_n \vdash B$$

**Theorem 2.10: The Curry-Howard Correspondence for  $\mathbf{TR}_0^{\text{bt}}$  and  $\mathbf{NJ}(\rightarrow)$ , Part 2.**

If  $\mathbf{NJ}(\rightarrow)$  proves

$$A_1, \dots, A_n \vdash B$$

Then for any distinct variables  $x_1, \dots, x_n$ , there is some bound-typed  $\lambda$ -term  $s$  such that  $\mathbf{TR}_0^{\text{bt}}$  proves

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

The proofs are trivial variants of Theorems 2.7 and 2.8, and are left to the reader.

Given the simplicity of  $\mathbf{TR}_0$  and  $\mathbf{NJ}(\rightarrow)$ , the Curry-Howard Correspondence and its applications given so far might seem somewhat superficial. But we will see over the next several chapters that the Curry-Howard Correspondence can be generalized and applied in further surprising ways. We have started this chapter with the most transparent form of the Curry-Howard Correspondence simply in order to have a firm and clear foundation on which to build more sophisticated versions.

### Exercises for Section 2.5

1. For each of the following sentences  $\tau$ , (i) construct a proof in  $\mathbf{NJ}(\rightarrow)$  of  $\tau$ , and (ii) using this proof and the method described in this section, generate a closed  $\lambda$ -term  $s$  and a proof in  $\mathbf{TR}_0$  of  $\vdash s : \tau$ . (In this way, you are showing that the type  $\tau$  is inhabited.)

(a)  $A \rightarrow (B \rightarrow B)$

(b)  $(A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C))$

(c)  $(A \rightarrow B) \rightarrow ((C \rightarrow A) \rightarrow (C \rightarrow B))$

2. Show that the sentence  $X \rightarrow (X \rightarrow X)$  is provable in  $\mathbf{NJ}(\rightarrow)$ , and that there exist terms  $s_1$  and  $s_2$  of the untyped  $\lambda$ -calculus such that both

$$\vdash s_1 : X \rightarrow (X \rightarrow X) \quad \text{and} \quad \vdash s_2 : X \rightarrow (X \rightarrow X)$$

are provable in  $\mathbf{TR}_0$ , but that  $s_1$  and  $s_2$  are not  $\beta\eta$ -equivalent. Thus, there is no sense in which the  $s$  postulated in Theorem 2.8 is unique.

## 2.6 Avoiding Misunderstandings: I

The Curry-Howard Correspondence tells us that a type  $A$  is inhabited (that is, there is a closed  $\lambda$ -term  $M$  which may be assigned type  $A$ ) iff  $A$  is a theorem of  $\mathbf{NJ}(\rightarrow)$ .

One might first find this result a little puzzling. Fix any types  $A$  and  $B$ . For example, suppose  $A$  is  $\mathbb{N}$  and  $B$  is  $\mathbb{B}$ . Surely I can come up with a computer program that, when given an input of type  $A$  returns an object of type  $B$ . For example, surely I can write a computer program that when given any natural number as input just returns the truth value  $\top$ . Is this not a computer program of type  $A \rightarrow B$ , for the values of  $A$  and  $B$  in question? And will this not ultimately mean that *any* type of the form  $A \rightarrow B$  is inhabited, regardless of whether  $A \rightarrow B$  is a theorem of  $\mathbf{NJ}(\rightarrow)$ ?

We must be careful here. We can certainly write down  $\lambda$ -terms that take as input an object of type  $A$  and that return an output of type  $B$ , for any  $A$  and  $B$ . The term  $\lambda x^A(y^B)$  is such a term. This term has type  $A \rightarrow B$ . This is true even if  $A \rightarrow B$  is not a theorem of  $\mathbf{NJ}(\rightarrow)$ . But this is not a counterexample to the Curry-Howard Correspondence, as the term  $\lambda x^A(y^B)$  is not closed, and so

$$\vdash \lambda x(y) : A \rightarrow B$$

is *not* provable in  $\mathbf{TR}_0$ . By contrast, you should be able to verify that

$$y : B \vdash \lambda x(y) : A \rightarrow B$$

is provable in  $\mathbf{TR}_0$ . But this poses no threat to the Curry-Howard Correspondence, as deleting  $\lambda$ -terms and keeping types gives us:

$$B \vdash A \rightarrow B$$

which is provable in  $\mathbf{NJ}(\rightarrow)$  for any  $A$  and  $B$ .

Alternatively, starting with the  $\lambda$ -term  $\lambda x^A(y^B)$ , we can bind the free variable  $y^B$  to obtain the closed  $\lambda$ -term  $\lambda y^B \lambda x^A(y^B)$ . But now this is a computer program with type  $B \rightarrow (A \rightarrow B)$ . This term *does* correspond to a theorem of intuitionistic logic, as you should be able to verify.

The lesson here is that although for arbitrary types  $A$  and  $B$  we can construct computer programs that take as inputs objects of type  $A$  and produce as output objects of type  $B$ , if the program involves free variables of any sort then it is *not* really a program of type  $A \rightarrow B$ , and so cannot count as refutation of the Curry-Howard Correspondence.

Still, one might have further worries. Let us go back to our computer program that takes any natural number as input, and returns the Boolean truth value  $\top$ . Or more generally, for any types  $A$  and  $B$ , fix some arbitrary object  $o^B$  of type  $B$ , and consider the  $\lambda$ -term  $\lambda x^A(o^B)$ . On any input of type  $A$ , this returns the fixed value  $o^B$ . The  $\lambda$ -term  $\lambda x^A(o^B)$  is *closed*. This closed  $\lambda$ -term looks like a solution to the inhabitation problem for  $A \rightarrow B$ . And so does this not mean after all that every type of the form  $A \rightarrow B$  is inhabited, regardless of whether  $A \rightarrow B$  is provable in  $\mathbf{NJ}(\rightarrow)$ ?

Again, once we think things through we see that there is no threat to the Curry-Howard Correspondence. The problem is that in the typed  $\lambda$ -calculus, we do *not* have constants  $o^B$  (or even  $\top^{\mathbb{B}}$ ) at our disposal. The term  $\lambda x^A(o^B)$  is *not* a term of the typed  $\lambda$ -calculus as we have defined it. Any constants we want to introduce must be explicitly defined. For some types, it is easy to define a constant of that type. But for other types we simply cannot.

For example, for any type of the form  $B \rightarrow B$  we can easily define an object  $o$  of type  $B \rightarrow B$ , just by regarding  $o^{B \rightarrow B}$  as an abbreviation for the closed  $\lambda$ -term  $\lambda x^B(x^B)$ . Understanding  $o^B$  as an abbreviation in this way, the term  $\lambda x^A(o^B)$  is perfectly acceptable. But this term has type  $A \rightarrow (B \rightarrow B)$ , which does indeed correspond to a theorem of  $\mathbf{NJ}(\rightarrow)$ .

By contrast, consider some simple type  $B$  such as  $\mathbb{N}$  or  $\mathbb{B}$ . Given only the machinery in the typed  $\lambda$ -calculus as we have defined it, we *cannot* define a closed term with type  $B$ . Any definition of any object in the typed  $\lambda$ -calculus by a closed term must have the form  $\lambda x^\rho(U)$  for for some  $U$ , and will thus have a type of the form  $\rho \rightarrow \rho'$  for some  $\rho$  and  $\rho'$ . We thus cannot define *any* specific object of *any* fundamental type with a closed term. So although we can use variables such as  $x^B$  in our programming language, we do not have prior access to any constants of this type, nor can we define them. Of course, as we saw earlier we can define objects that behave like Boolean truth values or natural numbers, but these are actually objects of more complex types that merely mimic the Booleans or natural numbers. They do not provide us with constants of the form  $o^B$  for any simple type  $B$ . Thus, we cannot view  $\lambda x^A(o^B)$  as a  $\lambda$ -term that solves the inhabitation problem for  $A \rightarrow B$  for any simple type  $B$ .

In fact, it is a consequence of the Curry-Howard Correspondence that we can define a closed term of type  $B$  iff  $B$  is a theorem of  $\mathbf{NJ}(\rightarrow)$ . (This further shows how restricted a programming language the typed  $\lambda$ -calculus is.) Thus, we can use terms of the form  $\lambda x^A o^B$  to solve the inhabitation problem for  $A \rightarrow B$  iff  $B$  is already a theorem of  $\mathbf{NJ}(\rightarrow)$ . But if  $B$  is already a theorem of  $\mathbf{NJ}(\rightarrow)$ , then so is  $A \rightarrow B$ . So again, there is no threat to the Curry-Howard Correspondence here.

## 2.7 More Consequences: Proof Normalization.

One important area in which the Curry-Howard correspondence can be put to very good use is *proof theory*. This is a rich and complex area of mathematics on its own, and we merely introduce the reader to a few important results without any claim to exhaustiveness. Some of the ideas introduced in this section will also be central in our later vindication of the **BHK** conditions.

Consider the following proof of  $Y \rightarrow (X \rightarrow Y)$  in  $\mathbf{NJ}(\rightarrow)$

$$\frac{\frac{\frac{X, X, Y \vdash Y}{X, Y \vdash X \rightarrow Y} (\rightarrow I) \quad X, Y \vdash X}{X, Y \vdash Y} (\rightarrow E) \quad \frac{X, Y \vdash Y}{Y \vdash X \rightarrow Y} (\rightarrow I)}{\vdash Y \rightarrow (X \rightarrow Y)} (\rightarrow I) \quad (38)$$

This proof is completely correct. However, it is perverse and inefficient. For example, the conclusion of the  $(\rightarrow E)$  inference is the sequent  $X, Y \vdash Y$ , which is a logical axiom with which we could have just begun the proof. There are other curious features of this proof that we will also turn to shortly.

Our goal is to discuss methods of proof simplification that could be used to at least somewhat simplify proofs like (38). There are presumably many different approaches that could be taken here, but as our goal is to develop very general methods of proof simplification, we want to focus on general principles rather than ad-hoc methods that only apply in very specific cases. We will see that there are in fact very general methods of proof simplification that can be applied not just to (38), but to any proof.

Let us then focus on a more subtle oddity of (38). Note that the proof begins with an application of the  $(\rightarrow I)$  rule introducing a conditional  $X \rightarrow Y$ , which is then immediately followed by an application of the  $(\rightarrow E)$  rule eliminating this very conditional. The result is the logical axiom  $X, Y \vdash Y$  with which we began the proof. There is something odd about introducing a conditional and then immediately eliminating it, and one might wonder whether this is the main cause of the inefficiency of (38).

Even more specifically, the proof (38) involves an application of the  $(\rightarrow E)$  rule whose major premise involves a conditional  $X \rightarrow Y$  that has just been introduced with the use of the  $(\rightarrow I)$  rule. Let us consider this pattern of reasoning more generally. We begin with the definition

**Definition 2.11: Normal Proof of  $\mathbf{NJ}(\rightarrow)$ .**

A proof in  $\mathbf{NJ}(\rightarrow)$  is *normal* just in case it contains no occurrence of a  $(\rightarrow\text{I})$  rule introducing a conditional  $X \rightarrow Y$  which is immediately followed by an application of the  $(\rightarrow\text{E})$  rule with this conditional  $X \rightarrow Y$  as its major premise.

The proof (38) is not normal. We would like to consider systematic ways of transforming it into a normal proof.

To this end, consider the following typing derivation that corresponds to (38) using the Curry-Howard correspondence

$$\frac{\frac{\frac{u : X, v : X, w : Y \vdash w : Y}{u : X, w : Y \vdash \lambda v(w) : X \rightarrow Y} \quad u : X, w : Y \vdash u : X}{u : X, w : Y \vdash \lambda v(w)(u) : Y}}{w : Y \vdash \lambda u(\lambda v(w)(u)) : X \rightarrow Y}}{\vdash \lambda w \lambda u(\lambda v(w)(u)) : Y \rightarrow (X \rightarrow Y)} \quad (39)$$

Just as there is something inefficient and unsatisfactory about (38), one would expect there to be something inefficient and unsatisfactory about (39). And indeed, there is something unsatisfactory about (39) that is very easy to identify. The typing derivation (39) shows that the type  $Y \rightarrow (X \rightarrow Y)$  is inhabited. The example it provides of a term that can be given this type is  $\lambda w \lambda u(\lambda v(w)(u))$ . But this term is not in  $\beta$ -normal form and can be simplified. When it is  $\beta$ -reduced, we get the  $\beta$ -normal form term  $\lambda w \lambda u(w)$ . Theorem 1.44 tells us that whenever an untyped  $\lambda$ -term  $M$  can be assigned a type and  $M \rightarrow_{\beta} N$ , then  $N$  can also be assigned that same type. Thus, we could have much more easily have shown that the type  $Y \rightarrow (X \rightarrow Y)$  was inhabited by showing that the term  $\lambda w \lambda u(w)$  can be given this type. In fact, such a typing derivation can easily be constructed

$$\frac{\frac{w : Y, u : X \vdash w : Y}{w : Y \vdash \lambda u(w) : X \rightarrow Y}}{\vdash \lambda w \lambda u(w) : Y \rightarrow (X \rightarrow Y)} \quad (40)$$

Removing the terms from this typing derivation and keeping the types then gives the following  $\mathbf{NJ}(\rightarrow)$  proof of  $Y \rightarrow (X \rightarrow Y)$

$$\frac{\frac{Y, X \vdash Y}{Y \vdash X \rightarrow Y} (\rightarrow\text{I})}{\vdash Y \rightarrow (X \rightarrow Y)} (\rightarrow\text{I}) \quad (41)$$

This is a very much simpler proof of  $Y \rightarrow (X \rightarrow Y)$ . Note moreover that (41) is a normal proof.

In fact, there is a very general method of transforming non-normal proofs to normal proofs lurking here. Before stating the main theorem that shows the general phenomenon, we need a lemma



**Lemma 2.12: The Subterm Lemma for  $\mathbf{TR}_0$** 

In any typing derivation of  $\mathbf{TR}_0$  with conclusion of the form  $\Gamma \vdash s : A$ , if  $s'$  is a subterm of  $s$  then the typing derivation contains a sequent of the form  $\Gamma' \vdash s' : B$ , where  $\Gamma' \subseteq \Gamma$ .

The proof of this lemma is a straightforward induction on the construction of typing derivations  $\Gamma \vdash s : A$  and is left to the exercises.

Armed with this lemma, the following theorem can be easily proven

**Theorem 2.13**

Given a typing derivation  $T$  in  $\mathbf{TR}_0$  whose conclusion is a sequent of the form  $\Delta \vdash g : C$ , let  $P$  be the  $\mathbf{NJ}(\rightarrow)$  proof obtained from  $T$  by deleting all terms and keeping types. Then  $P$  is a normal proof iff  $g$  is in  $\beta$ -normal form.

This theorem shows that there is an intimate connection between non-normal proofs and typing derivations of terms not in  $\beta$ -normal form. The proof is straightforward

**Proof**

First we show that if  $g$  is not in  $\beta$ -normal form, then  $P$  is not normal. Suppose then that  $T$  is a typing derivation with conclusion  $\Delta \vdash g : C$ , where  $g$  is not in  $\beta$ -normal form. It follows from Lemma 2.12 that  $T$  contains a sequent of the form

$$\Gamma \vdash \lambda u(e)(f) : D.$$

The step prior to this sequent must be the (App) rule

$$\frac{\begin{array}{c} \vdots \\ \Gamma \vdash \lambda u(e) : E \rightarrow D \end{array} \quad \begin{array}{c} \vdots \\ \Gamma \vdash f : E \end{array}}{\Gamma \vdash \lambda u(e)(f) : D} \text{ (App)}$$

and the step prior to the leftmost premise of this inference must be the (Abs) rule

$$\frac{\begin{array}{c} \vdots \\ \Gamma, u : E \vdash e : D \end{array}}{\Gamma \vdash \lambda u(e) : E \rightarrow D} \text{ (Abs)} \quad \begin{array}{c} \vdots \\ \Gamma \vdash f : E \end{array}}{\Gamma \vdash \lambda u(e)(f) : D} \text{ (App)}$$

Deleting all terms and keeping types then gives the  $\mathbf{NJ}(\rightarrow)$  proof  $P$

$$\frac{\frac{\vdots}{\bar{\Gamma}, E \vdash D} (\rightarrow\text{I}) \quad \frac{\vdots}{\bar{\Gamma} \vdash E} (\rightarrow\text{E})}{\bar{\Gamma} \vdash D} (\rightarrow\text{E})$$

for some list  $\bar{\Gamma}$  of formulae. This proof  $P$  is not normal.

Next we show that if  $P$  is not normal, then  $g$  is not in  $\beta$ -normal form. Suppose then that  $P$  has the form

$$\frac{\frac{\frac{\vdots}{\Gamma, A \vdash B} (\rightarrow\text{I}) \quad \frac{\vdots}{\Gamma \vdash A} (\rightarrow\text{E})}{\Gamma \vdash B} (\rightarrow\text{E})}{\vdots} \Delta \vdash C$$

Any typing derivation  $T$  from which this is derived must have the form

$$\frac{\frac{\frac{\vdots}{\bar{\Gamma}, u : A \vdash e : B} \quad \frac{\vdots}{\bar{\Gamma} \vdash f : A}}{\bar{\Gamma} \vdash \lambda u(e) : A \rightarrow B} \quad \bar{\Gamma} \vdash f : A}{\bar{\Gamma} \vdash \lambda u(e)(f) : B} \quad \vdots}{\bar{\Delta} \vdash g : C}$$

where  $u$  is a variable,  $e$ ,  $f$  and  $g$  are  $\lambda$ -terms, and  $\bar{\Gamma}$  and  $\bar{\Delta}$  are contexts.

Because in any typing derivation, every  $\lambda$ -term that appears in the proof is  $\alpha$ -equivalent to a subexpression of any  $\lambda$ -term lower in the proof, it follows that  $g$  has  $\lambda u(e)(f)$  (or something  $\alpha$ -equivalent to it) as a subexpression. Thus  $g$  is not in  $\beta$ -normal form.

Theorem 2.13 gives us a method for transforming any non-normal proof into a normal proof, by a combination of the Curry-Howard Correspondence and  $\beta$ -reduction. For given any non-normal proof  $P$ , we may transform it into a typing derivation  $T$  that contains  $P$  as its ‘skeleton’. By Theorem 2.13, this will be a typing derivation for a term  $e$  that is not in  $\beta$ -normal form, and it will assign  $e$  some type  $\tau$ . If  $e$   $\beta$ -reduces to  $e'$  where  $e'$  is in  $\beta$ -normal form, then  $e'$  can also be assigned the type  $\tau$ . We can thus produce a typing derivation that shows that  $e'$  can be assigned the type  $\tau$ . By Theorem 2.13, the skeleton of this typing derivation will be a normal proof  $P'$  (with the same conclusion as  $P$ ). We have thus transformed the non-normal proof  $P$  into a normal proof  $P'$ . The procedure is completely mechanical and requires no ingenuity. This procedure is often called ‘proof normalization’. Precisely this procedure transforms the proof (38) into the simpler proof (39). In the exercises, you will apply it to some other non-normal proofs.

Note then that we have the following result:

**Theorem 2.14: Proof Normalization Theorem for  $\mathbf{NJ}(\rightarrow)$**

If  $A_1, \dots, A_n \vdash B$  has a proof in  $\mathbf{NJ}(\rightarrow)$ , then it has a normal proof in  $\mathbf{NJ}(\rightarrow)$ .

**Proof**

Suppose  $P_1$  is a non-normal proof of  $A_1, \dots, A_n \vdash B$ . By the Curry-Howard Correspondence (in particular, Theorem 2.8), this may be transformed into a typing derivation with conclusion

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

for some  $\lambda$ -term  $s$ . By Theorem 2.13,  $s$  is not in  $\beta$ -normal form. Let  $s'$  be the  $\beta$ -normal form of  $s$ . By Theorem 1.44 (the Subject Reduction Theorem) there is a typing derivation of

$$x_1 : A_1, \dots, x_n : A_n \vdash s' : B$$

in  $\mathbf{TR}_0$ . Because  $s'$  is in  $\beta$ -normal form, it follows from Theorem 2.13 that the  $\mathbf{NJ}(\rightarrow)$  proof obtained by taking this typing derivation and removing all terms and keeping types is a normal proof of  $A_1, \dots, A_n \vdash B$ .

Normal proofs have a simple structure that makes them very useful in proving facts about  $\mathbf{NJ}(\rightarrow)$  or  $\mathbf{TR}_0$ . To see how, we first need another definition. Proofs in  $\mathbf{NJ}(\rightarrow)$  are *trees*. Such trees consists of *branches*, where a branch is a linear sequence of sequents that extends from some leaf of the tree (i.e., a logical axiom), and is such that if it includes some sequent then it includes the sequent immediately below it. In this way, a branch extends from a leaf of the tree all the way to the very bottom of the tree (i.e., the conclusion of the proof.) Consider for example the proof:

$$\frac{\frac{A, A \rightarrow B \vdash A \quad A, A \rightarrow B \vdash A \rightarrow B}{A, A \rightarrow B \vdash B} (\rightarrow E)}{A \vdash (A \rightarrow B) \rightarrow B} (\rightarrow I) \quad (42)$$

This has exactly two branches:

$$\begin{array}{lll} \mathbf{Branch 1:} & A, A \rightarrow B \vdash A & A, A \rightarrow B \vdash B \quad A \vdash (A \rightarrow B) \rightarrow B \\ \mathbf{Branch 2:} & A, A \rightarrow B \vdash A \rightarrow B & A, A \rightarrow B \vdash B \quad A \vdash (A \rightarrow B) \rightarrow B \end{array}$$

The following definition is useful

**Definition 2.15: Major Branch**

A *major branch* of a proof is a branch of the proof with the property that if the conclusion of an inference lies on the branch and that inference has a major premise, then that major premise also lies on in the branch.

You should be able to verify that in the proof just given, **Branch 1** is not a major branch, and **Branch 2** is a major branch.

Each proof contains at least one major branch: simply start at the bottom of the tree and move upwards from the conclusion of rule to its premises, picking the major premise whenever one reaches a ( $\rightarrow$ E) rule and one has two premises to pick from. Proceeding upwards in this way, one will eventually reach a leaf of the tree. In this way, one traces out a major branch. In the case of NJ( $\rightarrow$ ) proofs, the major branch is unique, though this will not hold true in more complicated systems we consider later.

Note that in a normal proof, on a major branch the ( $\rightarrow$ I) rule cannot be immediately followed by the ( $\rightarrow$ E) rule, as this would introduce the pattern:

major branch

$$\frac{\frac{\frac{\vdots}{\Gamma, A \vdash B}}{\Gamma \vdash A \rightarrow B} (\rightarrow I) \quad \frac{\frac{\vdots}{\Gamma \vdash A}}{\Gamma \vdash A} (\rightarrow E)}{\Gamma \vdash B} (\rightarrow E)$$

which would mean that the original proof was non-normal. Thus, on a major branch the ( $\rightarrow$ I) rule can only be followed by the ( $\rightarrow$ I) rule. This means that a major branch can only consist of a sequence of ( $\rightarrow$ E) rules followed by a sequence of ( $\rightarrow$ I) rules. For example, in the normal proof (42), **Branch 2** consists of a single ( $\rightarrow$ E) rule followed by a single ( $\rightarrow$ I) rule.

We have thus proven

**Lemma 2.16: The Major Branch Lemma.**

In any normal proof in NJ( $\rightarrow$ ), a major branch will consist of a (possibly empty) sequence of ( $\rightarrow$ E) rules followed by a (possibly empty) sequence of ( $\rightarrow$ I) rules.

This modest looking lemma turns out to be very useful. For example, with it one is easily able to prove the following

**Theorem 2.17: Consistency of  $\mathbf{NJ}(\rightarrow)$ .**

No atomic sentence  $X$  is provable in  $\mathbf{NJ}(\rightarrow)$ .

The only atomic sentences in the language used in  $\mathbf{NJ}(\rightarrow)$  are the propositional variables  $X, Y, Z, \dots$  and the sentence  $\perp$ . This theorem tells us that these sentences are not theorems of  $\mathbf{NJ}(\rightarrow)$ . This theorem may thus be viewed as a type of consistency proof for  $\mathbf{NJ}(\rightarrow)$ .

**Proof**

Suppose  $\mathbf{NJ}(\rightarrow)$  proves  $S$  where  $S$  is either a propositional variable or the symbol  $\perp$ . We then know from Theorem 2.14 that there is a normal proof  $P$  whose conclusion is  $\vdash S$ . The last step of  $P$  cannot be the  $(\rightarrow I)$  rule, because  $S$  does not have the form of a conditional. So using Lemma 2.16, we have that any major branch of  $P$  must consist solely of a sequence of  $(\rightarrow E)$  rules. But the  $(\rightarrow E)$  rule leaves the left hand side of a sequent unchanged. Because the left hand side of every leaf is non-empty, it follows that if a major branch of  $P$  consists only of  $(\rightarrow E)$  rules, then the left hand side of its bottommost sequent must also be non-empty. But the left hand side of  $\vdash S$  is empty, which is a contradiction.

Some further intriguing applications of the Major Branch Lemma will be given in the exercises.

The simplicity of  $\mathbf{NJ}(\rightarrow)$  might lead one to think that all this effort and terminology to prove these theorems is perhaps overkill. But will see that the basic ideas here can be generalized quite straightforwardly to help us analyze much stronger proof systems in a very similar way. As with the Curry-Howard Correspondence, we start with the simplest version of these ideas in order to have a firm and clear foundation on which to generalize them later.

**Exercises for Section 2.7**

1. Each of the following proofs in  $\mathbf{NJ}(\rightarrow)$  is *not* normal. Use the techniques of this section to transform them into normal proofs. In particular, begin by decorating the proofs in such a way to produce proofs in  $\mathbf{TR}_0$  of  $\vdash M : \tau$ , where  $\tau$  is the conclusion of the original  $\mathbf{NJ}(\rightarrow)$  proof. You should notice that  $M$  is not in  $\beta$ -normal form. Reduce  $M$  to a  $\beta$ -normal form  $M'$ , and produce a derivation of  $\vdash M' : \tau$  in  $\mathbf{TR}_0$ . This later typing derivation should encode a normal proof of  $\tau$ , and you should verify that this is the case.

(a)

$$\frac{\frac{A, A \vdash A}{A \vdash A \rightarrow A} (\rightarrow I) \quad A \vdash A}{\vdash A \rightarrow A} (\rightarrow E)$$

(b)

$$\frac{\frac{A, A, B \vdash B}{A, B \vdash A \rightarrow B} (\rightarrow I) \quad A, B \vdash A}{A, B \vdash B} (\rightarrow E) \quad \frac{A, B \vdash B}{B \vdash A \rightarrow B} (\rightarrow I) \quad \frac{B \vdash A \rightarrow B}{\vdash B \rightarrow (A \rightarrow B)} (\rightarrow I)$$

2. Show that the following are not derivable in  $\mathbf{NJ}(\rightarrow)$ , where  $X$  and  $Y$  are distinct atomic sentences. (Hint: use the Major Branch Lemma.) In each case, you make invoke the result from earlier cases.

- (i)  $X, (X \rightarrow Y) \rightarrow X \vdash Y$
- (ii)  $(X \rightarrow Y) \rightarrow X \vdash X \rightarrow Y$
- (iii)  $\vdash ((X \rightarrow Y) \rightarrow X) \rightarrow X$  (Pierce's Rule).

(side note: Pierce's Rule is provable in classical logic.)

3. Using the Major Branch Lemma, show that every closed  $\lambda$ -term in  $\beta$ -normal form has the form

$$\lambda x_1 \lambda x_2 \dots \lambda x_n (y(E_1)(E_2) \dots (E_m))$$

where  $n \geq 1$ ,  $y$  is a variable,  $m \geq 0$ , and each  $E_1, \dots, E_m$  is in  $\beta$ -normal form.

4. We have seen that a proof of  $\Delta \vdash E : X$  in  $\mathbf{TR}_0$  in which the  $\lambda$ -term  $E$  is in  $\beta$ -normal form corresponds to a proof in  $\mathbf{NJ}(\rightarrow)$  in which a certain pattern of reasoning does not appear. What additional pattern of reasoning does not appear if the  $\lambda$ -term  $E$  is assumed to be in  $\beta\eta$ -normal form?

5. Prove Lemma 2.12. (Hint: proceed by induction on the construction of typing derivations  $\Gamma \vdash s : A$ .)

6. Is there a normal proof of  $\mathbf{NJ}(\rightarrow)$  that includes a sequent  $\Gamma, X \vdash X$  not as a leaf of the proof?

7. Is there a normal proof of  $\mathbf{NJ}(\rightarrow)$  in which an occurrence of the  $(\rightarrow I)$  rule is immediately followed by an occurrence of the  $(\rightarrow E)$  rule, where it is not necessarily required that the conclusion of the  $(\rightarrow I)$  rule be the major premise of the  $(\rightarrow E)$  rule?

## 2.8 Adding Types: Ordered Pairs

The Curry-Howard correspondence shows that there is a perhaps surprising connection between certain natural concepts in the  $\lambda$ -calculus and the proof theory of  $\mathbf{NJ}(\rightarrow)$ . The deeper significance of the Curry-Howard correspondence however comes from the fact that these ideas can be generalized so that they hold in logical systems far stronger than  $\mathbf{NJ}(\rightarrow)$ . We begin this process of generalization now.

We have already argued that the typed  $\lambda$ -calculus can be viewed as a kind of (typed) program-

ming language. It is natural to want to expand the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  and its corresponding system  $\mathbf{TR}_0$  of typing rules to include more complex types that repeatedly arise in computer programming.

In a computer program, it is common to want to store two distinct items into a single data structure – an ordered pair – for subsequent manipulation. To capture this practice, given two untyped  $\lambda$ -terms  $x$  and  $y$  (i.e., two ‘computer programs’), we allow the formation of a new  $\lambda$ -term (i.e., a new computer program)  $\langle x, y \rangle$ , understood as the ordered pair with  $x$  as its first element and  $y$  as its second element. In order to extract information from ordered pairs, it is natural to also introduce into our programming language the functions  $\pi_1$  and  $\pi_2$  that take an ordered pair to its first and second elements respectively. Thus, in addition to our old formation rules, we allow the following new formation rules for  $\lambda$ -terms:

1. If  $x$  and  $y$  are  $\lambda$ -terms, then so is  $\langle x, y \rangle$ .
2. If  $x$  is a  $\lambda$ -term, then so are  $\pi_1(x)$  and  $\pi_2(x)$ .

So for example, in this expanded version of the  $\lambda$ -calculus we have terms like the following:

$$\begin{aligned} &\langle z, x \rangle \\ &\pi_1(\langle y, z \rangle) \\ &\langle z, x \rangle(\pi_1(\langle y, z \rangle)) \\ &\lambda x (\langle z, x \rangle(\pi_1(\langle y, z \rangle))) \\ &\pi_2(\lambda x (\langle z, x \rangle(\pi_1(\langle y, z \rangle)))) \end{aligned}$$

As before, we will omit parentheses where this causes no confusion, subject to our usual conventions. For example, the last term on our list can perhaps be written more simply as

$$\pi_2(\lambda x (\langle z, x \rangle \pi_1 \langle y, z \rangle)). \quad (43)$$

The free variables of  $\langle x, y \rangle$  are just all the free variables of  $x$  and all the free variables of  $y$ , and the free variables of  $\pi_1(x)$  and  $\pi_2(x)$  are just the free variables of  $x$ . So for example, in the  $\lambda$ -term (43) the free variables are  $z$  and  $y$ . If a variable occurs both free and bound, as in the term  $\langle x, \lambda x(x) \rangle$ , renaming of bound variables can be useful in determining which occurrences of which variables are free or bound. So we can rewrite this term in an  $\alpha$ -equivalent form as  $\langle x, \lambda y(y) \rangle$ , making it easier to see that the variable  $x$  is free and the variable  $y$  bound.

We expand the notion of  $\beta$ -reduction in this new programming language to allow the following reductions:

1.  $\pi_1(\langle M, N \rangle) \rightarrow_\beta M$
2.  $\pi_2(\langle M, N \rangle) \rightarrow_\beta N$

So for example, consider the term

$$\pi_1(\lambda x(x)(\langle y, z \rangle)).$$

This  $\beta$ -reduces to  $\pi_1(\langle y, z \rangle)$ , which then  $\beta$ -reduces to  $y$ . On the other hand, the term

$$\lambda x(x)(\pi_2(\langle y, z \rangle))$$

can be  $\beta$ -reduced in several ways. We have the following two  $\beta$ -reduction sequences:

$$\lambda x(x)(\pi_2(\langle y, z \rangle)) \rightarrow_{\beta} \lambda x(x)(z) \rightarrow_{\beta} z$$

$$\lambda x(x)(\pi_2(\langle y, z \rangle)) \rightarrow_{\beta} \pi_2(\langle y, z \rangle) \rightarrow_{\beta} z$$

Note that both  $\beta$ -reduction sequences give the same  $\beta$ -normal form term  $z$ . We will see later that the Church-Rosser Theorem holds for this new system, and so this convergence of results is not surprising.

Thus far, we have simply added some machinery to the untyped  $\lambda$ -calculus. The main task of the typed  $\lambda$ -calculus is the typing of untyped terms, and so in order to update the typed  $\lambda$ -calculus to accommodate our new constructions, we must think about how we should assign a type to  $\langle x, y \rangle$ . Here, there is a very natural idea: assuming that the type  $A$  can be assigned to  $x$ , and the type  $B$  can be assigned to  $y$ , we create a brand new type - the *product type*  $A \times B$  - and assert that it can be assigned to  $\langle x, y \rangle$ . Given two types  $A$  and  $B$ , the product type  $A \times B$  is simply the type of ordered pairs consisting of an element of  $A$  followed by an element of type  $B$ . This provides us with a brand new class of data structure. Up to now, the only way to form complex types from simple types was to build function types  $A \rightarrow B$ . Product types now give us a new way of forming more complex types.

The product type is governed by the following typing rules, which must be added to the rules of  $\mathbf{TR}_0$ :

$$\frac{\Gamma \vdash x : A \quad \Gamma \vdash y : B}{\Gamma \vdash \langle x, y \rangle : A \times B}$$

$$\frac{\Gamma \vdash x : A \times B}{\Gamma \vdash \pi_1(x) : A} \quad \frac{\Gamma \vdash x : A \times B}{\Gamma \vdash \pi_2(x) : B}$$

These rules tell us that if  $x$  has type  $A$  and  $y$  has type  $B$ , then the ordered pair  $\langle x, y \rangle$  has type  $A \times B$ . Conversely, if  $x$  has type  $A \times B$ , then  $\pi_1(x)$  has type  $A$  and  $\pi_2(x)$  has type  $B$ . Nothing should be surprising about these rules.

We have expanded the set of  $\lambda$ -terms, as well as the set of possible types a  $\lambda$ -term can be assigned. As a result, old terms can be assigned new types. For example, the  $\lambda$ -term  $\lambda x(x)$  can now be assigned the type  $(A \times B) \rightarrow (A \times B)$  for any types  $A$  and  $B$ , as the following typing derivation shows

$$\frac{x : A \times B \vdash x : A \times B}{\vdash \lambda x(x) : (A \times B) \rightarrow (A \times B)}$$

New terms can also be assigned new types. For example, the closed  $\lambda$ -term

$$\langle \lambda x(x), \lambda y(y) \rangle$$

can be given the type  $(A \rightarrow A) \times (B \rightarrow B)$  for any  $A$  and  $B$ , as the following typing derivation shows:



$$\frac{\frac{x : A \vdash x : A}{\vdash \lambda x(x) : A \rightarrow A} \quad \frac{y : B \vdash y : B}{\vdash \lambda y(y) : B \rightarrow B}}{\vdash \langle \lambda x(x), \lambda y(y) \rangle : (A \rightarrow A) \times (B \rightarrow B)}$$

For a slightly more complex example, consider the closed  $\lambda$ -term

$$\lambda x(\pi_1 \langle \lambda y(x), x \rangle).$$

This term can be given type  $A \rightarrow (B \rightarrow A)$  for any  $A$  and  $B$ , as the following typing derivation shows:

$$\frac{\frac{\frac{x : A, y : B \vdash x : A}{x : A \vdash \lambda y(x) : B \rightarrow A} \quad x : A \vdash x : A}{x : A \vdash \langle \lambda y(x), x \rangle : (B \rightarrow A) \times A}}{x : A \vdash \pi_1 \langle \lambda y(x), x \rangle : B \rightarrow A}}{\vdash \lambda x(\pi_1 \langle \lambda y(x), x \rangle) : A \rightarrow (B \rightarrow A)}$$

As before,  $\lambda$ -terms with free variables (i.e., open  $\lambda$ -terms) can still be typed, but the relevant typing derivation will have a non-empty context in its conclusion, indicating the types of the free variables. So for example, we might expect that  $\pi_1 \langle \lambda x(x), z \rangle$  can be assigned the type  $A \rightarrow A$  for any  $A$ . The following typing derivation confirms this:

$$\frac{\frac{z : B, x : A \vdash x : A}{z : B \vdash \lambda x(x) : A \rightarrow A} \quad z : B \vdash z : B}{z : B \vdash \langle \lambda x(x), z \rangle : (A \rightarrow A) \times B}}{z : B \vdash \pi_1 \langle \lambda x(x), z \rangle : A \rightarrow A}$$

Note that the conclusion of this typing derivation contains the non-empty context  $z : B$  declaring the type of the free variable  $z$ .

### Exercises for Section 2.8

1. Write each of the following  $\lambda$ -terms in  $\beta$ -normal form.

- (a)  $\pi_1(\pi_2(\langle x, y \rangle))$
- (b)  $\lambda x(\pi_1(\pi_1(\langle x, y \rangle)))(\pi_2(z))$
- (c)  $\lambda x(\pi_2(\pi_1(\langle x, y \rangle)))(\langle y, z \rangle)$

2. Each of the following untyped  $\lambda$ -terms can be assigned types. Find types for them, and confirm your result by producing an appropriate typing derivation.

- (a)  $\lambda x(\pi_1(\langle x, x \rangle))$
- (b)  $\lambda x(\pi_1(\langle y, x \rangle))$  (note: this is an open  $\lambda$ -term)
- (c)  $\lambda y(\langle \lambda x(\pi_1(\langle x, y \rangle)), y \rangle)$

## 2.9 Do We Need Product Types?

As an aside that may be skipped if the reader wishes, one might wonder whether we really need to add new machinery to the  $\lambda$ -calculus in order to encode ordered pairs. For example, given the untyped  $\lambda$ -terms  $E$  and  $F$ , consider the untyped  $\lambda$ -term

$$\lambda x(x(E)(F)).$$

Call this term **PAIR**( $E, F$ ). We will see that this  $\lambda$ -term ‘stores’  $E$  and  $F$ . In particular, letting **FIRST** be the  $\lambda$ -term  $\lambda y(y(\lambda u(\lambda v(u))))$  and letting **SECOND** be the  $\lambda$ -term  $\lambda y(y(\lambda u(\lambda v(v))))$ , and assuming  $x, v$  are not free in  $E$  or  $F$ , we have

$$\begin{aligned} \mathbf{FIRST}(\mathbf{PAIR}(E, F)) &= \lambda y(y(\lambda u(\lambda v(u))))(\lambda x(x(E)(F))) \\ &\rightarrow_{\beta} \lambda x(x(E)(F))(\lambda u(\lambda v(u))) \\ &\rightarrow_{\beta} \lambda u(\lambda v(u))(E)(F) \\ &\rightarrow_{\beta} \lambda v(E)(F) \\ &\rightarrow_{\beta} E \end{aligned}$$

and

$$\begin{aligned} \mathbf{SECOND}(\mathbf{PAIR}(E, F)) &= \lambda y(y(\lambda u(\lambda v(v))))(\lambda x(x(E)(F))) \\ &\rightarrow_{\beta} \lambda x(x(E)(F))(\lambda u(\lambda v(v))) \\ &\rightarrow_{\beta} \lambda u(\lambda v(v))(E)(F) \\ &\rightarrow_{\beta} \lambda v(v)(F) \\ &\rightarrow_{\beta} F \end{aligned}$$

Thus, **FIRST** and **SECOND** ‘extract’  $E$  and  $F$  respectively from **PAIR**( $E, F$ ). So **PAIR**( $E, F$ ) looks like it can play the role of an ordered pair, and **FIRST** and **SECOND** can play the role of  $\pi_1$  and  $\pi_2$ .

This is fine as far as the untyped  $\lambda$ -calculus goes – if the untyped  $\lambda$ -calculus were all we cared about, we could indeed dispense with product types in precisely this way. But unfortunately, this approach fails for the typed  $\lambda$ -calculus, as **FIRST**, **SECOND** and **PAIR**( $E, F$ ) cannot be typed in a consistent manner.

To see why, suppose  $E$  has type  $A$  and  $F$  has type  $B$ . What type will **PAIR**( $E, F$ ) - i.e.,  $\lambda x(x(E)(F))$  receive? For the subexpression  $x(E)(F)$  to be grammatical,  $x$  must have type  $A \rightarrow (B \rightarrow C)$  for some type  $C$ . Then  $x(E)(F)$  will have type  $C$ , and  $\lambda x(x(E)(F))$  will have type  $(A \rightarrow (B \rightarrow C)) \rightarrow C$ . So **PAIR**( $E, F$ ) will have type  $(A \rightarrow (B \rightarrow C)) \rightarrow C$  for some  $C$ .

Because **FIRST**(**PAIR**( $E, F$ )) simplifies to  $E$  which is of type  $A$ , and **PAIR**( $E, F$ ) has type  $(A \rightarrow (B \rightarrow C)) \rightarrow C$ , it follows that **FIRST** must have type

$$((A \rightarrow (B \rightarrow C)) \rightarrow C) \rightarrow A.$$

Writing out **FIRST** as the  $\lambda$ -term  $\lambda y(y(\lambda u\lambda v(u)))$ , it follows that this term must be typed in the following way:

$$\lambda y^{(A \rightarrow (B \rightarrow C)) \rightarrow C} (y^{(A \rightarrow (B \rightarrow C)) \rightarrow C} (\lambda u\lambda v(u)))^A \quad (44)$$

(Here we use the fact that if  $\lambda s(t)$  has type  $E \rightarrow F$ , then  $s$  must have type  $E$  and  $t$  must have type  $F$ .) How can we assign types to  $u$  and  $v$  for the expression (44) to make sense? In order for the subexpression

$$y^{(A \rightarrow (B \rightarrow C)) \rightarrow C} (\lambda u\lambda v(u)) \quad (45)$$

to be grammatical,  $\lambda u\lambda v(u)$  must have type  $A \rightarrow (B \rightarrow C)$ , and so then expression (45) gets type  $C$ . But looking at (44), the subexpression (45) was supposed to get type  $A$ . This means that  $C$ , which up to now was unspecified, must be  $A$ , and so **PAIR**( $E, F$ ) has type  $(A \rightarrow (B \rightarrow A)) \rightarrow A$ .

Because **SECOND**(**PAIR**( $E, F$ )) simplifies to  $F$  which is of type  $B$ , and **PAIR**( $E, F$ ) has type  $(A \rightarrow (B \rightarrow A)) \rightarrow A$ , it follows that **SECOND** must have type

$$((A \rightarrow (B \rightarrow A)) \rightarrow A) \rightarrow B.$$

Writing out **SECOND** as the  $\lambda$ -term  $\lambda y(y(\lambda u\lambda v(v)))$ , it follows that this term must be typed in the following way:

$$\lambda y^{(A \rightarrow (B \rightarrow A)) \rightarrow A} (y^{(A \rightarrow (B \rightarrow A)) \rightarrow A} (\lambda u\lambda v(v)))^B \quad (46)$$

This means that the subexpression  $\lambda u\lambda v(v)$  of (46) must get the type  $A \rightarrow (B \rightarrow A)$ . This can only happen as follows:

$$\lambda u^A \lambda v^B (v^A)$$

but then  $v$  is assigned both  $A$  and  $B$  as types, which need not be equal.

So we cannot type the terms **FIRST**, **SECOND** and **PAIR** in a consistent way. This problem plagues all attempts to try to simulate ordered pairs and projection functions within the  $\lambda$ -calculus  $\lambda_0$ . (You will see the deeper reason for this in the exercises.) It is for this reason that we need to introduce pairing and projections into the  $\lambda$ -calculus directly if we hope to incorporate this sort of machinery into a version of the *typed*  $\lambda$ -calculus.

### Exercises for Section 2.9

1. Show that there do not exist typable  $\lambda_0$ -terms **FIRST**, **SECOND** and **PAIR**( $x, y$ ) such that the following conditions all hold:

- (i) **FIRST** and **SECOND** are closed  $\lambda_0$ -terms
- (ii) **PAIR**( $x, y$ ) is an open  $\lambda_0$ -term with free variables  $x$  and  $y$ ,
- (iii) **FIRST**(**PAIR**( $x, y$ ))  $\rightarrow_\beta x$
- (iv) **SECOND**(**PAIR**( $x, y$ ))  $\rightarrow_\beta y$

(Hint: suppose to the contrary that such typable  $\lambda_0$ -terms exist. Then assuming that the variables  $x$  and  $y$  have types  $A$  and  $B$  respectively, we will be able to assign **PAIR**( $x, y$ ) some type  $P$  (which may depend on  $A$  and  $B$ .) Because we can input **PAIR**( $x, y$ ) to **FIRST** and

the result is of type  $A$ , it follows that **FIRST** must in fact have the type  $P \rightarrow A$ . Likewise, **SECOND** must in fact have the type  $P \rightarrow B$ . Thus, the following three sequents must be derivable in  $\mathbf{TR}_0$ :

$$x : A, y : B \vdash \mathbf{PAIR}(x, y) : P \quad \vdash \mathbf{FIRST} : P \rightarrow A \quad \vdash \mathbf{SECOND} : P \rightarrow B$$

By the Curry-Howard Correspondence, the following three sequents must then also be derivable in  $\mathbf{NJ}(\rightarrow)$ :

$$A, B \vdash P \quad \vdash P \rightarrow A \quad \vdash P \rightarrow B$$

Assuming without loss of generality that  $A, B$  are distinct propositional variables, show that no sentence  $P$  exists with this property. This in effect amounts to showing that conjunction is not definable from implication in  $\mathbf{NJ}(\rightarrow)$ . The Major Branch Lemma applied to  $\vdash P \rightarrow A$  and  $\vdash P \rightarrow B$  will help here. In particular, using the Major Branch Lemma you should be able to show that  $P$  has both the forms

$$C_n \rightarrow \dots \rightarrow (C_2 \rightarrow (C_1 \rightarrow B))$$

and

$$D_m \rightarrow \dots \rightarrow (D_2 \rightarrow (D_1 \rightarrow A)),$$

which is a contradiction.)

## 2.10 Adding Types: Sums

From the point of view of a computer programmer, there are other types it is natural to want to manipulate beyond function types and product types. For example, given two types  $A$  and  $B$ , it is natural to want to consider something like their union. We will denote this  $A + B$ . It is convenient to think of this as not just a union, but as something a little more complex. In particular, unlike a plain union, it is convenient to think of each element of  $A + B$  as an object from  $A$  or  $B$  which also has attached to it a label indicating whether it was drawn from  $A$  or  $B$ . Just by looking at this label we can tell whether the object is from  $A$  or  $B$ , without looking at the object itself. In such a case, we call  $A + B$  a *sum type*.

One way (amongst many) to think of the elements of  $A + B$  is as something like ordered pairs  $\langle a, b \rangle$  in which  $a = 1$  and  $b$  is an element of  $A$ , or  $a = 2$  and  $b$  is an element of  $B$ . As a result, merely by looking at the label  $a$  we can know whether  $b$  came from  $A$  or  $B$ . On this way of doing things, given an object  $x$  in  $A$ ,  $x$  itself is not an element of  $A + B$ , but  $\langle 1, x \rangle$  is. Likewise, given an object  $x$  in  $B$ ,  $x$  itself is not an element of  $A + B$ , but  $\langle 2, x \rangle$  is. This is one way of implementing the idea of ‘labels’ being attached to terms in  $A + B$ . So for example, on this way of doing things the sum type  $A + A$  consists of the objects of the form  $\langle i, x \rangle$  where  $i$  is 1 or 2, and  $x$  an object of type  $A$ .

We will *not* proceed in this way, but will instead think of these ‘labels’ in slightly more abstract terms. We do this by assuming that in our programming language we have a type of **if-then-else**

command that, when input with an element  $x$  of type  $A + B$ , does one thing if  $x$  is something with the label 1 attached to it, and something else if it sees that  $x$  is something with the label 2 attached to it. More specifically, we assume we have access to commands in our computer language that say something like:

If  $x$  is an object  $z$  with the label 1 attached, then run the program  $P[z/v]$ ; i.e., run the program  $P$  with the free variable  $v$  replaced with  $z$ , and if  $x$  is an object  $z$  with the label 2 attached, run the program  $Q[z/w]$ ; i.e., run the program  $Q$  with the free variable  $w$  replaced with  $z$ . (47)

For the sake of completeness, we can also specify some sort of behavior when  $x$  has neither the forms in question (e.g., we could have the program halt and return an error message, or just output some fixed value and continue), though in what follows the details of what happens in this case will not matter.

We denote the command described in (47) by

$$\mathbf{case}(x, [v]P, [w]Q).$$

Of course, the programs  $P$  and  $Q$  will here be  $\lambda$ -terms, and by ‘running’ a program, we just mean making the appropriate substitutions and then  $\beta$ -reducing (or  $\beta\eta$ -reducing) the resulting expression.

The notation  $[v]P$  is perhaps unfamiliar. This is really just the program (or  $\lambda$ -term)  $P$ , but with a particular variable  $v$  singled out as special. This special variable  $v$  is the one whose free occurrences we will replace with  $x$  should the operation described in (47) require us to. Likewise for  $[w]Q$ .

So for example, consider the command:

$$\mathbf{case}(x, [z]\lambda u(uz)[z], [v]v)$$

If  $x$  is the term  $\langle 1, \bar{x} \rangle$ , then this produces the result  $\lambda u(u\bar{x})$ . If  $x$  is the term  $\langle 2, \bar{x} \rangle$ , this produces the result  $\bar{x}$ .

In the case in which  $u$  does not occur free in  $P$  (for example,  $P$  is a closed term) and  $v$  does not occur free in  $Q$ , then the command

$$\mathbf{case}(x, [u]P, [v]Q)$$

just tells us to produce the output  $P$  if the label attached to  $x$  is 1, and the output  $Q$  if the label attached to  $x$  is 2. So for example, the command

$$\mathbf{case}(x, [v]\lambda u(u), [v]\lambda w(w))$$

outputs  $\lambda u(u)$  if the label attached to  $x$  is 1, and  $\lambda w(w)$  if the label attached to  $x$  is 2. This sort of command mimics the **if-then-else** command structures found in many programming languages.

Note that once we have access to these sorts of **case** commands, we have no reason to care about the precise mechanics by which the ‘label’ 1 or 2 is attached to an element of a sum type

$A + B$ . As long as we have the ability to execute one of a given pair of commands in the case that the label is 1 and the other in case the label is 2, we can tell what label occurs in  $x$  simply by looking at the outcome of  $\mathbf{case}(x, [u]P, [v]Q)$  for appropriately chosen  $P$  and  $Q$ . There is no need to talk about the mechanics of label attachment beyond that.

To see how this works in full detail, it will also be convenient to introduce an operator  $\mathbf{in}_1$  that takes a  $\lambda$ -term  $x$  to  $x$  with the label 1 attached, as well as an operator  $\mathbf{in}_2$  that takes the  $\lambda$ -term  $x$  to  $x$  with the label 2 attached. Again, the precise mechanism by which labels are being attached will not matter. Although we can think of  $\mathbf{in}_1(x)$  as  $\langle 1, x \rangle$  and  $\mathbf{in}_2(x)$  as  $\langle 2, x \rangle$ , it is more clean to simply regard  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$  as basic expressions that attach labels to  $x$ , without the need to specify the precise mechanics by which they do so.

In addition to **1.** and **2.** above, we then have the following new formation rules for untyped  $\lambda$ -terms:

- 3.** If  $x$  is a  $\lambda$ -term, then  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$  are  $\lambda$ -terms.
- 4.** If  $x, P$  and  $Q$  are  $\lambda$ -terms and  $u, v$  variables, then  $\mathbf{case}(x, [u]P, [v]Q)$  is a  $\lambda$ -term.

So for example, in this expanded version of the  $\lambda$ -calculus we have terms like the following:

$$\begin{aligned} & \mathbf{in}_1(z) \\ & \langle \mathbf{in}_1(z), x \rangle \\ & \langle \mathbf{in}_1(z), x \rangle \pi_1(y) \\ & \lambda x (\langle \mathbf{in}_1 z, x \rangle \pi_1(y)) \\ & \mathbf{case}(\lambda x (\langle \mathbf{in}_1 z, x \rangle \pi_1(y)), [u]\pi_2(u), [w]\langle w, v \rangle) \end{aligned}$$

As always, we omit parentheses where no confusion results.

The free variables of  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$  are just the free variables of  $x$ . The free variables of  $\mathbf{case}(x, [u]P, [v]Q)$  are all the free variables of  $x$ , all the free variables of  $P$  except for  $u$ , and all the free variables of  $Q$  except for  $v$ . So in the expression

$$\mathbf{case}(\lambda x (\langle \mathbf{in}_1 z, x \rangle \pi_1(y)), [u]\pi_2(u), [w]\langle w, v \rangle)$$

you should be able to verify that the free variables are  $z, y$  and  $v$ .

Note in particular that any free occurrences of  $v$  in  $P$  become bound in the expression  $[v]P$ . In this way,  $[v]$  is a ‘binding operator’ that, much like  $\lambda$ -abstraction, takes a free variable and renders it bound.

The obvious  $\beta$ -reduction rules to add at this point are

- 3.**  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q) \rightarrow_{\beta} P[R/u]$
- 4.**  $\mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q) \rightarrow_{\beta} Q[R/v]$

So for example, the terms

$$\mathbf{case}(\lambda x (\mathbf{in}_1(x))(y), [z]z, [y]\lambda z(y))$$

and

$$\lambda x (\mathbf{case}(x, [u]\pi_2(u), [v]\pi_1(v))) \mathbf{in}_2(\langle y, z \rangle)$$

may be  $\beta$ -reduced as follows:

$$\begin{aligned} \mathbf{case}(\lambda x(\mathbf{in}_1(x))(y), [z]z, [y]\lambda z(y)) &\rightarrow_{\beta} \mathbf{case}(\mathbf{in}_1(y), [z]z, [y]\lambda z(y)) \\ &\rightarrow_{\beta} z[\mathbf{in}_1(y)/z] \\ &= \mathbf{in}_1(y) \end{aligned}$$

$$\begin{aligned} \lambda x(\mathbf{case}(x, [u]\pi_2(u), [v]\pi_1(v))) \mathbf{in}_2(\langle y, z \rangle) &\rightarrow_{\beta} \mathbf{case}(\mathbf{in}_2(\langle y, z \rangle), [u]\pi_2(u), [v]\pi_1(v)) \\ &\rightarrow_{\beta} \pi_1(v)[\langle y, z \rangle/v] \\ &= \pi_1(\langle y, z \rangle) \\ &\rightarrow_{\beta} y \end{aligned}$$

Next, let us think about what sorts of types we should assign  $\lambda$ -terms of the form  $\mathbf{in}_1(x)$ ,  $\mathbf{in}_2(x)$ , and  $\mathbf{case}(x, P[v_1], Q[v_2])$ . We begin by thinking about the typing rules for  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$ . Given types  $\sigma$  and  $\tau$  and an object  $x$  of type  $\sigma$ , we use  $\mathbf{in}_1(x)$  to attach a label to  $x$  indicating that it has type  $\tau$ , thereby creating an object of the sum type  $\sigma + \tau$ . With this in mind, one might suspect that the following are the typing rules for  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$ :

$$\frac{\Gamma \vdash x : \sigma}{\Gamma \vdash \mathbf{in}_1(x) : \sigma + \tau} \qquad \frac{\Gamma \vdash x : \tau}{\Gamma \vdash \mathbf{in}_2(x) : \sigma + \tau} \quad (48)$$

However, these rules cannot be quite right. Focusing on the leftmost rule, given an  $x$  of type  $\sigma$ , for *which* type  $\tau$  is  $\mathbf{in}_1(x)$  of type  $\sigma + \tau$ ? The type  $\tau$  simply appears out of nowhere in the conclusion of this rule. So if  $x$  is of type  $\mathbb{N}$ , is  $\mathbf{in}_1(x)$  of type  $\mathbb{N} + \mathbb{N}$ , or  $\mathbb{N} + \mathbb{B}$ , or  $\mathbb{N} + (\mathbb{N} \rightarrow \mathbb{N})$ ? We want every introduced expression to have a definite type. There must therefore be a determinate answer to the question of what type  $\mathbf{in}_1(x)$  has.

To resolve this ambiguity, in the context of typed computations we add superscripts to the notation  $\mathbf{in}_1(x)$ , and write  $\mathbf{in}_1^{\sigma, \tau}$  to indicate that this is an operation that takes an object of type  $\sigma$  and produces an output of type  $\sigma + \tau$ . Likewise,  $\mathbf{in}_2^{\sigma, \tau}$  is an operation that takes an object of type  $\tau$  and produces an output of type  $\sigma + \tau$ . In the context of the untyped  $\lambda$ -calculus, we do not need these superscripts, and can take at face value the rule that if  $x$  is a  $\lambda$ -term, then  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$  are  $\lambda$ -terms. In the context of the typed  $\lambda$ -calculus, however, we must regard not only all variables as typed, but also all expressions of the form  $\mathbf{in}_1$  and  $\mathbf{in}_2$  as typed. So we will *not* regard the term  $\mathbf{in}_1(x^{\mathbb{N}})$  as fully typed. By contrast, the term  $\mathbf{in}_1^{\mathbb{N}, \mathbb{B}}(x^{\mathbb{N}})$  will be regarded as fully typed, and moreover has the type  $\mathbb{N} + \mathbb{B}$ .

The correct typing rules for  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$  are thus:

$$\frac{\Gamma \vdash x : \sigma}{\Gamma \vdash \mathbf{in}_1^{\sigma, \tau}(x) : \sigma + \tau} \qquad \frac{\Gamma \vdash x : \tau}{\Gamma \vdash \mathbf{in}_2^{\sigma, \tau}(x) : \sigma + \tau} \quad (49)$$

Having said all this, in the context of a typing derivation the superscripts on  $\mathbf{in}_1$  and  $\mathbf{in}_2$  can typically be determined from context. So for example, in the leftmost of the rules (48), the only

superscripts that could be attached to  $\mathbf{in}_1$  to render the inference grammatical would be the superscripts  $\sigma, \tau$ . When the superscripts on  $\mathbf{in}_1$  or  $\mathbf{in}_2$  can be uniquely determined by context, we typically omit them for the sake of compact notation. In such cases, the superscripts should still be understood as present, but merely ‘invisible’. In cases in which there is any ambiguity or possibility of confusion, we include all superscripts. With these conventions in mind, our original rules (48) are perfectly acceptable; but without these conventions in mind, they do not make sense.

We now turn to the question of what type we should assign  $\mathbf{case}(x, [u]P, [v]Q)$ . The problem is that in general this term can get two different types. Suppose, for example, that whenever  $z$  has type  $\sigma$ ,  $P[z/u]$  has type  $\rho_1$ . And suppose that whenever  $z$  has type  $\tau$ ,  $Q[z/v]$  has type  $\rho_2$ . Now suppose that  $x$  is an element of type  $\sigma + \tau$ . Then depending on whether  $x$  is an object  $z$  of type  $\sigma$  with the label 1 attached to it or an object  $z$  of type  $\tau$  with the label 2 attached to it,  $\mathbf{case}(x, [u]P, [v]Q)$  will either have type  $\rho_1$  or type  $\rho_2$ . And so it looks like we cannot make a definitive statement in advance about which type  $\mathbf{case}(x, [u]P, [v]Q)$  will receive.

Suppose however that  $\rho_1$  and  $\rho_2$  (i.e., the type of  $P[z/u]$  if  $z$  has type  $\sigma$ , and the type of  $Q[z/v]$  if  $z$  has type  $\tau$ ) are actually the same type  $\rho$ . In this case, we *can* definitively say that  $\mathbf{case}(x, [u]P, [v]Q)$  gets the type  $\rho$ . This suggests the following typing rule

$$\frac{\Gamma \vdash x : \sigma + \tau \quad \Gamma, u : \sigma \vdash P : \rho \quad \Gamma, v : \tau \vdash Q : \rho}{\Gamma \vdash \mathbf{case}(x, [u]P, [v]Q) : \rho} \quad (50)$$

The rules (49) and (50) capture the defining properties of sum types. Note that they do so without presupposing anything about the mechanism by which labels are attached to elements in sum types. From a type theoretic point of view, directs sum types are then typically thought of as defined by these typing rules, rather than in terms of ordered pairs.

As before, armed with these rules we can type terms of the untyped  $\lambda$ -calculus. For example, consider the term

$$\lambda x(\mathbf{case}(x, [y]\mathbf{in}_1(y), [z]\mathbf{in}_2(z)))$$

This can be assigned the type  $(A+B) \rightarrow (A+B)$  for any  $A$  and  $B$ , as the following typing derivation shows:

$$\frac{x : A+B \vdash x : A+B \quad \frac{x : A+B, y : A \vdash y : A}{x : A+B, y : A \vdash \mathbf{in}_1(y) : A+B} \quad \frac{x : A+B, z : B \vdash z : B}{x : A+B, z : A \vdash \mathbf{in}_2(z) : A+B}}{x : A+B \vdash \mathbf{case}(x, [y]\mathbf{in}_1(y), [z]\mathbf{in}_2(z)) : A+B} \quad \frac{}{\vdash \lambda x(\mathbf{case}(x, [y]\mathbf{in}_1(y), [z]\mathbf{in}_2(z))) : (A+B) \rightarrow (A+B)}$$

In this typing derivation,  $\mathbf{in}_1$  should of course be understood as  $\mathbf{in}_1^{A,B}$  throughout, and  $\mathbf{in}_2$  should be understood as  $\mathbf{in}_2^{A,B}$  throughout. Consider next the open  $\lambda$ -term

$$\mathbf{case}(x, [y]\langle y, u \rangle, [z]\langle u, z \rangle)$$

This  $\lambda$ -term has  $x$  and  $u$  as free variables. If  $x$  has type  $A+A$  and  $u$  has type  $A$ , then this  $\lambda$ -term can be assigned type  $A \times A$ , as the following typing derivation sketch shows:



$$\frac{x : A+A, u : A \vdash x : A+A \quad \begin{array}{c} \vdots \\ x : A+A, u : A, y : A \vdash \langle y, u \rangle : A \times A \end{array} \quad \begin{array}{c} \vdots \\ x : A+A, u : A, z : A \vdash \langle u, z \rangle : A \times A \end{array}}{x : A+A, u : A \vdash \mathbf{case}(x, [y]\langle y, u \rangle, [z]\langle u, z \rangle) : A \times A}$$

You should easily be able to reconstruct the omitted upper parts  $\vdots$  of this typing derivation. Note also that the conclusion of this typing derivation contains the context  $x : A+A, u : A$  declaring the type of the free variables  $x$  and  $u$  of  $\mathbf{case}(x, [y]\langle y, u \rangle, [z]\langle u, z \rangle)$ .

### Exercises for Section 2.10

1. Write each of the following  $\lambda$ -terms in  $\beta$ -normal form.

- (a)  $\mathbf{case}(\mathbf{in}_1(x), \lambda z(y)[y], z[z])$
- (b)  $\lambda v(\mathbf{case}(v, [y]y, [z]\lambda x(z)(x)))(\mathbf{in}_2(x))$
- (c)  $\lambda v(\mathbf{case}(v(z), [z]\pi_1(z), [z]\pi_2(z)))(\lambda x(\mathbf{in}_2(\langle x, y \rangle)))$

2. Each of the following untyped  $\lambda$ -terms can be assigned types. Find types for them, and confirm your result by producing an appropriate typing derivation. If the  $\lambda$ -term is not closed, the conclusion of your typing derivation will of course have a non-empty context.

- (a)  $\lambda x(\mathbf{case}(x, [y]\pi_1(y), [z]z))$
- (b)  $\mathbf{case}(x, [y]\langle y, u \rangle, [z]\langle v, z \rangle)$
- (c)  $\mathbf{case}(x, [u]u(y), [v]z(v))$

## 2.11 Adding Types: The Empty Type

In type theory, it is also useful to sometimes be able to refer to the *empty* type, which we denote  $0$ . The idea is that no expression should have this type. We express this somewhat modestly with the idea that if a term has type  $0$ , then from this term we can construct an object of any type we please. In particular, given an  $x$  of type  $0$  and a type  $A$ , we can construct a term of type  $A$ , which we call  $\mathbf{abs}^A(x)$  (the ‘absurd’ object of type  $A$ .) We thus add the following rule for  $\lambda$ -term formation in the untyped  $\lambda$ -calculus:

5. If  $x$  is a  $\lambda$ -term, then so is  $\mathbf{abs}(x)$ .

and the typing rule:

$$\frac{\Gamma \vdash x : 0}{\Gamma \vdash \mathbf{abs}^A(x) : A}$$

Because  $\lambda$ -terms must have determinate types in the typed  $\lambda$ -calculus, it is necessary to add a superscript  $A$  to  $\mathbf{abs}(x)$ , just as it was necessary to add superscripts to the  $\mathbf{in}$  operators. However, just as with the  $\mathbf{in}$  operators we will omit superscripts when they are obvious from the context. We may thus rewrite our rule as

$$\frac{\Gamma \vdash x : 0}{\Gamma \vdash \mathbf{abs}(x) : A}$$

The free variables of  $\mathbf{abs}(x)$  are just the free variables of  $x$ . There is no corresponding new rule of  $\beta$ -reduction to be introduced.

Using this expanded  $\lambda$ -calculus and new typing rule, we can for example construct a  $\lambda$ -term that can be given the type  $(A \& (A \rightarrow 0)) \rightarrow B$  for any types  $A$  and  $B$  as follows:

$$\frac{\frac{x : A \& (A \rightarrow 0) \vdash x : A \& (A \rightarrow 0)}{x : A \& (A \rightarrow 0) \vdash \pi_2(x) : A \rightarrow 0} \quad \frac{x : A \& (A \rightarrow 0) \vdash x : A \& (A \rightarrow 0)}{x : A \& (A \rightarrow 0) \vdash \pi_1(x) : A}}{x : A \& (A \rightarrow 0) \vdash \pi_2(x)(\pi_1(x)) : 0} \quad \frac{x : A \& (A \rightarrow 0) \vdash \mathbf{abs}(\pi_2(x)(\pi_1(x))) : B}{\vdash \lambda x(\mathbf{abs}(\pi_2(x)(\pi_1(x)))) : (A \& (A \rightarrow 0)) \rightarrow B}$$

In this derivation,  $\mathbf{abs}$  is of course to be understood as  $\mathbf{abs}^B$ .

### Exercises for Section 2.11

1. Each of the following untyped  $\lambda$ -terms can be assigned types. Find types for them, and confirm your result by producing an appropriate typing derivation.

(a)  $\lambda x(\mathbf{abs}(\pi_1(x)))$

(b)  $\lambda x(\mathbf{abs}(\mathbf{case}(x, y[y], \pi_1(z)[z])))$

## 2.12 $\eta$ Rules for Product and Sum Types.

In addition to the new rules of  $\beta$ -reduction we have introduced for product and sum types, there are also new rules of  $\eta$ -reduction appropriate for product and sum types. We introduce these in this section.

The primary way of constructing functions in the  $\lambda$ -calculus is by  $\lambda$ -abstraction over a  $\lambda$ -term. So given a  $\lambda$ -term  $t$ , we can abstract over it to obtain a function  $\lambda x(t)$ . But consider an object  $f$  of some function type  $A \rightarrow B$ . Is there anything that guarantees that  $f$  is obtained by abstracting over some  $\lambda$ -term? Identifying  $f$  with  $\lambda x(f(x))$  (where  $x$  is a variable not free in  $f$ ) tells us that  $f$  is indeed an abstraction over the term  $f(x)$  of type  $B$  with respect to a variable  $x$  of type  $A$ . So the rule of  $\eta$ -reduction

$$\lambda x(t(x)) \rightarrow_{\eta} t$$

(where  $x$  is not free in  $t$ ) tells us that *all* functions in the  $\lambda$ -calculus are abstractions over terms. More precisely, it allows us to say that every object  $f$  of type  $A \rightarrow B$  is  $\eta$ -equivalent to a  $\lambda$ -abstraction over a  $\lambda$ -term – that is,  $f$  is  $\eta$ -equivalent to  $\lambda x(f(x))$  (where we select a variable  $x$  that is not free in  $f$ .)

Likewise, the primary way of constructing objects of a product type  $A \times B$  is by taking an object  $t_1$  of type  $A$  and an object  $t_2$  of type  $B$  and forming the object  $\langle t_1, t_2 \rangle$ . But is there anything

that guarantees that any object  $t$  of type  $A \times B$  has this form? It is natural to identify any  $t$  of type  $A \times B$  with  $\langle \pi_1(t), \pi_2(t) \rangle$ . We therefore introduce the new rule of  $\eta$ -reduction

$$\langle \pi_1(t), \pi_2(t) \rangle \rightarrow_{\eta} t.$$

This rule guarantees that every object of type  $A \times B$  can be written in the form  $\langle t_1, t_2 \rangle$ . More precisely, it allows us to say that every object  $t$  of type  $A \times B$  is  $\eta$ -equivalent to an explicit ordered pair – namely,  $\langle \pi_1(t), \pi_2(t) \rangle$ .

Turning to sum types, the primary way of constructing objects of a sum type  $A + B$  is either by taking an object  $x$  of type  $A$  and forming  $\mathbf{in}_1(x)$ , or by taking an object  $y$  of type  $B$  and forming  $\mathbf{in}_2(y)$ . (This case is trickier than the previous two in that we have multiple options here.) But is there anything that guarantees that any object  $t$  of type  $A + B$  must have one of these two forms? Note that for any  $t$  of type  $A + B$ , it is natural to identify  $t$  with

$$\mathbf{case}(t, [u]\mathbf{in}_1(u), [u]\mathbf{in}_2(u)). \quad (51)$$

To see why, we evaluate the term (51), considering each case individually. First, if  $t$  is of the form  $\mathbf{in}_1(z)$ , then the whole term simplifies to  $\mathbf{in}_1(u)[z/u]$ , which is  $\mathbf{in}_1(z)$ , and so the final result is  $t$  itself. Likewise, if  $t$  is of the form  $\mathbf{in}_2(z)$ , then the whole term simplifies to  $\mathbf{in}_2(u)[z/u]$ , which is  $\mathbf{in}_2(z)$ , and so the final result is again  $t$  itself. So either way, the term (51) simplifies to  $t$ .

Because a term of the form  $\mathbf{case}(t, [u]s_1(u), [v]s_2(v))$  has the form  $s_1(z)$  or  $s_2(z)$  for some  $z$  (assuming that  $u$  does not appear free in  $s_1$  or  $s_2$ ), identifying  $t$  with (51) then tells us that  $t$  has either the form  $\mathbf{in}_1(z)$  or  $\mathbf{in}_2(z)$  for some  $z$ , as desired. It is thus natural to introduce the new rule of  $\eta$ -reduction

$$\mathbf{case}(t, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) \rightarrow_{\eta} t. \quad (52)$$

This rule guarantees that every object of type  $A + B$  has either the form  $\mathbf{in}_1(t)$  or  $\mathbf{in}_2(t)$ .

Our  $\eta$ -rules tell us that objects of complex types have *canonical forms*: every element of a function type  $A \rightarrow B$  has the form  $\lambda x(t)$ , every element of a product type  $A \times B$  has the form  $\langle t_1, t_2 \rangle$ , and every element of a sum type  $A + B$  has the form  $\mathbf{in}_1(t)$  or  $\mathbf{in}_2(t)$ . Later on when we have a notion of equality in our theory we will return to the question of the meaning of the  $\eta$ -rules, but the present considerations will suffice for now.

One might be tempted to generalize (52) somewhat, and write

$$\mathbf{case}(t, [u]s[\mathbf{in}_1(u)/z], [v]s[\mathbf{in}_2(v)/z]) \rightarrow_{\eta} s[t/z] \quad (53)$$

where  $s, t$  are  $\lambda$ -terms and  $u, v, z$  variables such that  $u$  and  $v$  are not free in  $s$ . Note that (52) may be obtained from (53) by letting  $s$  be the term  $z$ . This rule looks entirely reasonable. Moreover, it has the virtue of allowing us to see a sense in which  $\lambda$ -terms such as

$$R(\mathbf{case}(z, [x]P, [y]Q)) \quad \text{and} \quad \mathbf{case}(z, [x]R(P), [y]R(Q))$$

are equivalent (where  $P, Q$  and  $R$  are  $\lambda$ -terms and  $x, y, z$  are variables such that neither  $x$  nor  $y$  are free in  $R$  and  $z$  is not free in  $P$  or  $Q$ .) That these are equivalent in some sense is intuitive

– it should not matter whether we divide into cases, apply  $P$  or  $Q$ , and then input the result into  $R$ , or whether we divide into cases, then apply  $R(P)$  or  $R(Q)$ . To see this more formally, substitute  $R(\mathbf{case}(z, [x]P, [y]Q))$  for  $s$  in (53), where  $u$  and  $v$  are assumed to not appear free in  $R(\mathbf{case}(z, [x]P, [y]Q))$ . This gives the expression

$$\mathbf{case}(t, [u]R(\mathbf{case}(z, [x]P, [y]Q))[in_1(u)/z], [v]R(\mathbf{case}(z, [x]P, [y]Q))[in_2(v)/z]) \quad (54)$$

Now, because  $z$  does not appear free in  $P$  or  $Q$ , the expression  $R(\mathbf{case}(z, [x]P, [y]Q))[in_1(u)/z]$  is  $R(\mathbf{case}(in_1(u), [x]P, [y]Q))$ , which  $\beta$ -reduces to  $R(P[u/x])$ . Because  $x$  is not free in  $R$ , this is  $R(P)[u/x]$ . Likewise,  $s[in_2(v)/z]$   $\beta$ -reduces to  $R(Q)[v/y]$ . Thus (54)  $\beta$ -reduces to

$$\mathbf{case}(z, [u]R(P)[u/x], [v]R(Q)[v/y]).$$

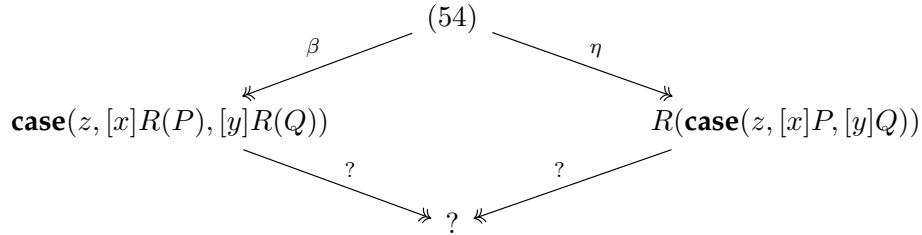
Renaming bound variables, we can rewrite this more simply as

$$\mathbf{case}(z, [x]R(P), [y]R(Q)). \quad (55)$$

But on the other hand, using the rule (53), the expression (54)  $\eta$ -reduces to  $R(\mathbf{case}(z, [x]P, [y]Q))$ . Because the same expression both  $\beta$ -reduces to  $\mathbf{case}(z, [x]R(P), [y]R(Q))$  and also  $\eta$ -reduces to  $R(\mathbf{case}(z, [x]P, [y]Q))$ , we then have

$$R(\mathbf{case}(z, [x]P, [y]Q)) \equiv_{\beta\eta} \mathbf{case}(z, [x]R(P), [y]R(Q)).$$

On one had, this is a pleasing result. On the other hand, without further modifications this leads to a violation of the Church-Rosser property for  $\beta\eta$ -reductions, as there is no obvious way to complete the ‘diamond’



While different approaches are possible to rectify this, the approach we will take is to prohibit the  $\eta$ -reduction rule (53). In some texts, rules such as (53) as re-classified as rules of ‘commutation’, along with reduction rules like

$$\begin{aligned} \mathbf{case}(t, [u]P, [v]Q)(R) &\rightarrow \mathbf{case}(t, [u]P(R), [v]Q(R)) \\ \lambda x(\mathbf{case}(t, [u]P, [v]Q)) &\rightarrow \mathbf{case}(t, [u]\lambda x(P), [v]\lambda x(Q)) \end{aligned}$$

and so on (with appropriate restrictions on free variables.) In general, such rules allow us to move certain expressions in and out of binding operators. Many tricky issues arise here, which we do not pursue further. For our purposes, we allow only the  $\eta$ -reduction rules:

$$\lambda x(t(x)) \rightarrow_{\eta} t$$

$$\langle \pi_1(t), \pi_2(t) \rangle \rightarrow_\eta t$$

$$\mathbf{case}(t, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) \rightarrow_\eta t$$

and no others. More formally, we have the following definitions

**Definition 2.18:  $\eta$ -reduction,  $\beta\eta$ -reduction, and  $\beta\eta$ -equivalence**

- (a)  $A \rightarrow_\eta B$  iff for some expressions  $E$  and  $F$
- (i)  $A$  has the form  $E \lambda x(t(x)) F$  and  $B$  has the form  $E t F$  (where  $x$  is not free in  $t$ ), or
  - (ii)  $A$  has the form  $E \langle \pi_1(t), \pi_2(t) \rangle F$  and  $B$  has the form  $E t F$ , or
  - (iii)  $A$  has the form  $E \mathbf{case}(t, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) F$  and  $B$  has the form  $E t F$ .
- (b)  $A \rightarrow_{\beta\eta} B$  iff  $A \rightarrow_\beta B$  or  $A \rightarrow_\eta B$ .
- (c)  $A \twoheadrightarrow_{\beta\eta} B$  (i.e.,  $A$   $\beta\eta$ -reduces to  $B$ ) iff for some sequence  $X_0, X_1, \dots, X_n$  of  $\lambda$ -terms we have that  $A = X_0$ ,  $B = X_n$ , and  $X_0 \rightarrow_{\beta\eta} X_1 \dots \rightarrow_{\beta\eta} X_n$ . This includes the degenerate case in which  $A$  and  $B$  are identical.
- (d)  $X$  and  $Y$  are  $\beta\eta$ -equivalent (in symbols,  $X =_{\beta\eta} Y$ ) iff one of the following holds:
- (i)  $X$  and  $Y$  are the same term (up to  $\alpha$ -equivalence), or
  - (ii) for some sequence of terms  $X_0, X_1, \dots, X_{n-1}, X_n$  with  $U = X_0$  and  $V = X_n$ ,

$$X_0 \sim_{\beta\eta} X_1, X_1 \sim_{\beta\eta} X_2, \dots, X_{n-1} \sim_{\beta\eta} X_n.$$

## 2.13 Collecting Things Together

We have added much machinery to the untyped  $\lambda$ -calculus  $\lambda_0$ , creating a more sophisticated untyped calculus with many rules. For convenience, we present these rules together.

**Definition 2.19: Rules of the untyped  $\lambda$ -calculus  $\lambda_1$ .**

Rules of  $\lambda$ -term formation

- (i) Any variable  $x, y, z, \dots$  is a  $\lambda$ -term.
- (ii) If  $M$  and  $N$  are  $\lambda$ -terms, then so is  $M(N)$ .
- (iii) If  $M$  is a  $\lambda$ -term and  $x$  a variable, then  $\lambda x(M)$  is a  $\lambda$ -term.
- (iv) If  $x$  and  $y$  are  $\lambda$ -terms, then so is  $\langle x, y \rangle$ .
- (v) If  $x$  is a  $\lambda$ -term, then so are  $\pi_1(x)$  and  $\pi_2(x)$ .
- (vi) If  $x$  is a  $\lambda$ -term, then so are  $\mathbf{in}_1(x)$  and  $\mathbf{in}_2(x)$ .
- (vii) If  $R, P$  and  $Q$  are  $\lambda$ -terms and  $u, v$  variables, then  $\mathbf{case}(R, [u]P, [v]Q)$  is a  $\lambda$ -term.

(viii) If  $x$  is a  $\lambda$ -term, then so is  $\mathbf{abs}(x)$ .

Rules of  $\beta$ -reduction

- (i)  $\lambda x(M)(N) \rightarrow_{\beta} M[N/x]$
- (ii)  $\pi_1 \langle M, N \rangle \rightarrow_{\beta} M$
- (iii)  $\pi_2 \langle M, N \rangle \rightarrow_{\beta} N$
- (iv)  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q) \rightarrow_{\beta} P[R/u]$
- (v)  $\mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q) \rightarrow_{\beta} Q[R/v]$

Rules of  $\eta$ -reduction

- (i)  $\lambda x(M(x)) \rightarrow_{\eta} M$  where  $x$  is not free in  $M$
- (ii)  $\langle \pi_1(M), \pi_2(M) \rangle \rightarrow_{\eta} M$
- (iii)  $\mathbf{case}(M, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) \rightarrow_{\eta} M$

We call this expanded version of the untyped  $\lambda$ -calculus  $\lambda_1$ .

To the untyped  $\lambda$ -calculus  $\lambda_1$ , we append the following set of typing rules:

**Definition 2.20: Rules of the typed  $\lambda$ -calculus  $\mathbf{TR}_1$ .**

- (i)  $\Gamma, v : \tau \vdash v : \tau$
- (ii) 
$$\frac{\Gamma \vdash x : \tau \rightarrow \sigma \quad \Gamma \vdash y : \tau}{\Gamma \vdash x(y) : \sigma}$$
- (iii) 
$$\frac{\Gamma, v : \tau \vdash x : \sigma}{\Gamma \vdash \lambda v(x) : \tau \rightarrow \sigma}$$
- (iv) 
$$\frac{\Gamma \vdash x : \sigma \quad \Gamma \vdash y : \tau}{\Gamma \vdash \langle x, y \rangle : \sigma \times \tau}$$
- (v) 
$$\frac{\Gamma \vdash x : \sigma \times \tau}{\Gamma \vdash \pi_1(x) : \sigma} \quad \frac{\Gamma \vdash x : \sigma \times \tau}{\Gamma \vdash \pi_2(x) : \tau}$$
- (vi) 
$$\frac{\Gamma \vdash x : \sigma}{\Gamma \vdash \mathbf{in}_1(x) : \sigma + \tau} \quad \frac{\Gamma \vdash x : \tau}{\Gamma \vdash \mathbf{in}_2(x) : \sigma + \tau}$$
- (vii) 
$$\frac{\Gamma \vdash x : \sigma + \tau \quad \Gamma, u : \sigma \vdash y : \rho \quad \Gamma, v : \tau \vdash z : \rho}{\Gamma \vdash \mathbf{case}(x, [u]y, [v]z) : \rho}$$
- (viii) 
$$\frac{\Gamma \vdash x : 0}{\Gamma \vdash \mathbf{abs}(x) : \tau}$$

where  $v, v_1, v_2$  are variables,  $x, y, z$  are terms of  $\lambda_1$ , and  $\sigma, \tau, \rho$  are types.

This system of typing rules is sometimes called  $\lambda(I)$ . We will refer to it as  $\mathbf{TR}_1$ . It is an expansion of the smaller typing system  $\lambda(\rightarrow)$  (i.e.,  $\mathbf{TR}_0$ ) considered earlier.

As before, there is a bound-typed version of this set of typing rules

**Definition 2.21: Rules of the typed  $\lambda$ -calculus  $\text{TR}_1^{\text{bt}}$ .**

- (i)  $\Gamma, v : \tau \vdash v : \tau$
- (ii) 
$$\frac{\Gamma \vdash x : \tau \rightarrow \sigma \quad \Gamma \vdash y : \tau}{\Gamma \vdash x(y) : \sigma}$$
- (iii) 
$$\frac{\Gamma, v : \tau \vdash x : \sigma}{\Gamma \vdash \lambda v^\tau(x) : \tau \rightarrow \sigma}$$
- (iv) 
$$\frac{\Gamma \vdash x : \sigma \quad \Gamma \vdash y : \tau}{\Gamma \vdash \langle x, y \rangle : \sigma \times \tau}$$
- (v) 
$$\frac{\Gamma \vdash x : \sigma \times \tau}{\Gamma \vdash \pi_1(x) : \sigma} \quad \frac{\Gamma \vdash x : \sigma \times \tau}{\Gamma \vdash \pi_2(x) : \tau}$$
- (vi) 
$$\frac{\Gamma \vdash x : \sigma}{\Gamma \vdash \mathbf{in}_1(x) : \sigma + \tau} \quad \frac{\Gamma \vdash x : \tau}{\Gamma \vdash \mathbf{in}_2(x) : \sigma + \tau}$$
- (vii) 
$$\frac{\Gamma \vdash x : \sigma + \tau \quad \Gamma, u : \sigma \vdash y : \rho \quad \Gamma, v : \tau \vdash z : \rho}{\Gamma \vdash \mathbf{case}(x, [u^\sigma]y, [v^\tau]z) : \rho}$$
- (viii) 
$$\frac{\Gamma \vdash x : 0}{\Gamma \vdash \mathbf{abs}(x) : \tau}$$

where  $v, v_1, v_2$  are variables,  $x, y, z$  are terms of  $\lambda_1$ , and  $\sigma, \tau, \rho$  are types.

There is also a corresponding generalization of the Church-style  $\lambda$ -calculus:

**Definition 2.22: Rules of the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$ .**Types

The *propositional types* are those that may be obtained from the fundamental types  $\tau_1, \tau_2, \dots$  and the empty type 0 under the operations taking two types  $A$  and  $B$  to the types  $A \rightarrow B$ ,  $A \times B$  and  $A + B$ .

Terms

The set of *simply typed  $\lambda$ -terms* are defined inductively as follows:

- (i) For any variable  $v$  and any type  $\tau$ ,  $v^\tau$  is a  $\lambda$ -term of type  $\tau$ .
- (ii) If  $[M]^{\alpha \rightarrow \beta}$  and  $[N]^\alpha$  are  $\lambda$ -terms, then  $[[M]^{\alpha \rightarrow \beta}([N]^\alpha)]^\beta$  is a  $\lambda$ -term.
- (iii) If  $[M]^\beta$  is a  $\lambda$ -term and  $x^\alpha$  a variable, then  $[\lambda x^\alpha([M]^\beta)]^{\alpha \rightarrow \beta}$  is a  $\lambda$ -term.
- (iv) If  $[M]^\alpha, [N]^\beta$  are  $\lambda$ -terms, then  $[\langle [M]^\alpha, [N]^\beta \rangle]^{\alpha \times \beta}$  is a  $\lambda$ -term.
- (v) If  $[M]^{\alpha \times \beta}$  is a  $\lambda$ -term, then  $[\pi_1([M]^{\alpha \times \beta})]^\alpha$  and  $[\pi_2([M]^{\alpha \times \beta})]^\beta$  are  $\lambda$ -terms.
- (vi) If  $[M]^\alpha$  and  $[N]^\beta$  are  $\lambda$ -terms, then  $[\mathbf{in}_1^{\alpha, \beta}([M]^\alpha)]^{\alpha + \beta}$  and  $[\mathbf{in}_2^{\alpha, \beta}([N]^\beta)]^{\alpha + \beta}$  are  $\lambda$ -terms.
- (vii) If  $[R]^{\alpha + \beta}$ ,  $[P]^\gamma$ , and  $[Q]^\gamma$  are  $\lambda$ -terms and  $u^\alpha, v^\beta$  variables, then

- [**case**([ $R$ ] $^{\alpha+\beta}$ , [ $u^\alpha$ ] $[P]^\gamma$ , [ $v^\beta$ ] $[Q]^\gamma$ )] $^\gamma$  is a  $\lambda$ -term.  
 (viii) If [ $M$ ] $^0$  is a  $\lambda$ -term, then so is [**abs** $^\alpha$ ( $[M]^0$ )] $^\alpha$ .

#### Rules of $\beta$ -reduction

- (i)  $\lambda x^\sigma ([M]^\tau)([N]^\sigma) \rightarrow_\beta [M]^\tau [[N]^\sigma / x^\sigma]$
- (ii)  $\pi_1(\langle [M]^\sigma, [N]^\tau \rangle^{\sigma \times \tau}) \rightarrow_\beta [M]^\sigma$
- (iii)  $\pi_2(\langle [M]^\sigma, [N]^\tau \rangle^{\sigma \times \tau}) \rightarrow_\beta [N]^\tau$
- (iv) **case**(**in** $_{1}^{\alpha, \beta}$ ( $[S]^\alpha$ ) $^{\alpha+\beta}$ , [ $u^\alpha$ ] $[P]^\gamma$ , [ $v^\beta$ ] $[Q]^\gamma$ )  $\rightarrow_\beta [P]^\gamma [[S]^\alpha / u^\alpha]$
- (v) **case**(**in** $_{2}^{\alpha, \beta}$ ( $[S]^\beta$ ) $^{\alpha+\beta}$ , [ $u^\alpha$ ] $[P]^\gamma$ , [ $v^\beta$ ] $[Q]^\gamma$ )  $\rightarrow_\beta [Q]^\gamma [[S]^\beta / v^\beta]$

#### Rules of $\eta$ -reduction

- (i)  $\lambda x^\sigma ([M]^{\sigma \rightarrow \tau}(x^\sigma)) \rightarrow_\eta [M]^{\sigma \rightarrow \tau}$  where  $x^\sigma$  is not free in  $M$
- (ii)  $\langle \pi_1([M]^{\sigma \times \tau}), \pi_2([M]^{\sigma \times \tau}) \rangle \rightarrow_\eta [M]^{\sigma \times \tau}$
- (iii) **case**( $[M]^{\sigma+\tau}$ , [ $u^\sigma$ ]**in** $_{1}^{\sigma, \tau}$ ( $u^\sigma$ ), [ $v^\tau$ ]**in** $_{2}^{\sigma, \tau}$ ( $v^\tau$ ))  $\rightarrow_\eta [M]^{\sigma+\tau}$

We call this version of the Church-style typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$ . Although in this definition we have written everything out in full detail, as with  $\lambda_0^{\text{type}}$  we will in practice omit brackets and types for clarity when they can be easily reconstructed. In particular, as before it will suffice to indicate the types of bound variables when they are first introduced, as well as the types of free variables. Note however that in our expanded version of the  $\lambda$ -calculus we now have *two* sources of bound variables. First, there are variables connected with  $\lambda$ -operators, as before. Second, we have variables that are bound within a **case** term. Recall that in

$$\mathbf{case}(x, [u]P, [v]Q),$$

$[u]$  and  $[v]$  are regarded as ‘binding operators’ much like  $\lambda$ , and their scopes are all of  $P$  and  $Q$  respectively. Thus, if we are to type all bound variables at the moment the relevant binding operators are introduced, we must type the variables  $[u]$  and  $[v]$  in **case**( $x, [u]P, [v]Q$ ). For example, consider the open  $\lambda_1$ -term

$$\lambda y \mathbf{case}(x, [u]u, [v]vy)$$

We can type all the bound variables of this expression as follows:

$$\lambda y^B \mathbf{case}(x, [u^A]u, [v^{B \rightarrow A}]vy)$$

The free variable  $x$  obviously must have type  $A + (B \rightarrow A)$  (because in **case**( $x, [u]P, [v]Q$ ) the type of  $x$  must be the sum of the types of  $u$  and  $v$ ), and thus we may regard this term as a fully typed  $\lambda_1^{\text{type}}$ -term – explicitly writing down the type of the free variable  $x$  is in this case optional. However, where it does not produce excessive clutter, for maximum clarity we generally indicate the type of each free variable explicitly, even when it is otherwise deducible. So for the term in question we will generally write

$$\lambda y^B \mathbf{case}(x^{A+(B \rightarrow A)}, [u^A]u, [v^{B \rightarrow A}]vy)$$

That this is grammatically typed is shown by the following typing derivation of **TR** $_1$



$$\frac{\begin{array}{c} \vdots \\ x:A+(B \rightarrow A), y:B \vdash x:A+(B \rightarrow A) \quad x:A+(B \rightarrow A), y:B, u:A \vdash u:A \quad x:A+(B \rightarrow A), y:B, v:B \rightarrow A \vdash vy:A \end{array}}{\frac{x:A+(B \rightarrow A), y:B \vdash \mathbf{case}(x, [u]u, [v]vy):A}{x:A+(B \rightarrow A) \vdash \lambda y \mathbf{case}(x, [u]u, [v]vy):B \rightarrow A}}$$

or perhaps even more clearly by the following typing derivation of  $\mathbf{TR}_1^{\text{bt}}$

$$\frac{\begin{array}{c} \vdots \\ x:A+(B \rightarrow A), y:B \vdash x:A+(B \rightarrow A) \quad x:A+(B \rightarrow A), y:B, u:A \vdash u:A \quad x:A+(B \rightarrow A), y:B, v:B \rightarrow A \vdash vy:A \end{array}}{\frac{x:A+(B \rightarrow A), y:B \vdash \mathbf{case}(x, [u^A]u, [v^{B \rightarrow A}]vy):A}{x:A+(B \rightarrow A) \vdash \lambda y^B \mathbf{case}(x, [u^A]u, [v^{B \rightarrow A}]vy):B \rightarrow A}}$$

For another example, consider the similar untyped  $\lambda_1$ -term

$$\lambda y \mathbf{case}(xz, [u]u, [v]vy)$$

as before, the bound variables may be typed as follows

$$\lambda y^B \mathbf{case}(xz, [u^A]u, [v^{B \rightarrow A}]vy)$$

However, knowing that  $xz$  has type  $A + (B \rightarrow A)$  does not allow us to determine the types of the free variables  $x$  and  $z$ . In order for this term to be fully typed, we must therefore indicate the types of at least one of these terms explicitly. So for example, we regard the term

$$\lambda y^B \mathbf{case}(xz^D, [u^A]u, [v^{B \rightarrow A}]vy)$$

as a fully typed  $\lambda_1^{\text{type}}$ -term, as in this expression  $x$  must obviously have type  $D \rightarrow (A + (B \rightarrow A))$ . Just as before, for maximum clarity it is better to indicate the type of each free variable explicitly, in which case we write

$$\lambda y^B \mathbf{case}(x^{(A+(B \rightarrow A)) \rightarrow D} z^D, [u^A]u, [v^{B \rightarrow A}]vy).$$

The reader is invited to write out the corresponding typing derivation in either  $\mathbf{TR}_1$  or  $\mathbf{TR}_1^{\text{bt}}$ .

Note that in addition to familiar examples of untypable terms such as  $x(x)$ , in  $\lambda_1$  we now have a large class of new untypable terms. For example,

$$\pi_1(\lambda x(E))$$

cannot be typed for any  $\lambda$ -term  $E$ , because when fully typed,  $\pi_1$  only acts on objects whose types have the form  $A \times B$ , while  $\lambda x(E)$  always has a type of the form  $A \rightarrow B$ . In the exercises you will see examples of other terms in  $\lambda_1$  that cannot be typed for similar reasons.

There are many results about  $\lambda_1$ ,  $\lambda_1^{\text{type}}$  and  $\mathbf{TR}_1$  that are straightforward generalizations of earlier theorems about  $\lambda_0$ ,  $\lambda_0^{\text{type}}$  and  $\mathbf{TR}_0$ . In many cases, their proofs only involve only minor modifications of the corresponding proofs in the previous chapter. Occasionally, when a result about our expanded  $\lambda$ -calculus has a straightforward and essentially identical proof to that given

in the previous chapter, we simply use the result directly. (For example, the fact that every  $\lambda$ -term is  $\alpha$ -equivalent to one that obeys the Barendregt variable convention remains true in  $\lambda_1$ , and the proof is for all intents and purposes identical to that of the corresponding fact about  $\lambda_0$ .) Doing so helps us to avoid having to relist a mass of theorems every time we expand the  $\lambda$ -calculus. We state some of the more important and less trivial generalizations here, in each leaving the proof to the Appendix or the exercises.

**Lemma 2.23:**

For any  $\lambda_1$ -terms  $U, V, t$  and variable  $x$ , if  $U \rightarrow_\beta V$  then  $U[t/x] \rightarrow_\beta V[t/x]$ , and if  $U \rightarrow_\eta V$  then  $U[t/x] \rightarrow_\eta V[t/x]$ .

**Theorem 2.24: The Church-Rosser Theorem for  $\beta$ - and  $\beta\eta$ -reductions in  $\lambda_1$ .**

In the untyped  $\lambda$ -calculus  $\lambda_1$ , if  $X \rightarrow_\beta Y_1$  and  $X \rightarrow_\beta Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_\beta Z$  and  $Y_2 \rightarrow_\beta Z$ . Likewise, if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

**Corollary 2.25**

In  $\lambda_1$ ,  $U =_\beta V$  iff there is a term  $Z$  such that  $U \rightarrow_\beta Z$  and  $V \rightarrow_\beta Z$ , and  $U =_{\beta\eta} V$  iff there is a term  $Z$  such that  $U \rightarrow_{\beta\eta} Z$  and  $V \rightarrow_{\beta\eta} Z$ .

**Theorem 2.26: The Church-Rosser Theorem for  $\beta$ - and  $\beta\eta$ -reductions in  $\lambda_1^{\text{type}}$ .**

In the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$ , if  $X \rightarrow_\beta Y_1$  and  $X \rightarrow_\beta Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_\beta Z$  and  $Y_2 \rightarrow_\beta Z$ . Likewise, if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

**Corollary 2.27**

In  $\lambda_1^{\text{type}}$ ,  $U =_\beta V$  iff there is a term  $Z$  such that  $U \rightarrow_\beta Z$  and  $V \rightarrow_\beta Z$ , and  $U =_{\beta\eta} V$  iff there is a term  $Z$  such that  $U \rightarrow_{\beta\eta} Z$  and  $V \rightarrow_{\beta\eta} Z$ .

**Theorem 2.28: Strong Normalizability of  $\lambda_1^{\text{type}}$  under  $\beta$ -reduction.**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$$

Using  $\eta$ -reduction postponement (Lemma 2.34, to be discussed shortly), it can be shown as before that if a term has an infinite  $\beta\eta$ -reduction sequence, then it also has an infinite  $\beta$ -reduction sequence. (The analog for this in  $\lambda_0$  was proved in the exercises for section 1.7, and the argument goes through unchanged.) The following is therefore a consequence of Theorem 2.28

**Theorem 2.29: Strong Normalizability of  $\lambda_1^{\text{type}}$  under  $\beta\eta$ -reduction.**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta\eta} \tau_2 \rightarrow_{\beta\eta} \tau_3 \rightarrow_{\beta\eta} \dots$$

We also have

**Theorem 2.30**

Suppose  $E$  is a  $\lambda_1$ -term obeying the Barendregt variable convention, and that  $\Gamma \vdash E : X$  in  $\mathbf{TR}_1$ . Then there is a proof of  $\Gamma' \vdash E : X$  in  $\mathbf{TR}_1$  for some  $\Gamma' \subseteq \Gamma$  such that (i) there is no silent replacement of  $\lambda$ -terms by  $\alpha$ -equivalent terms in this proof, (ii) the variables whose types are declared in  $\Gamma'$  are precisely the free variables of  $E$ , and (iii) all the elements of  $\Gamma'$  are distinct, and (iv) every variable declaration in the proof is either an element of  $\Gamma'$ , or has the form  $v : \tau$  where  $v$  is a bound variable of  $E$ .

**Theorem 2.31: Soundness and Completeness of  $\mathbf{TR}_1$ .**

Suppose  $E$  is an untyped  $\lambda_1$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i)  $E$  can be typed into a  $\lambda_1^{\text{type}}$  term with overall type  $X$  in which the free variables  $v_1, \dots, v_n$  are assigned types  $\tau_1, \dots, \tau_n$  respectively.
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_1$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

**Theorem 2.32: Soundness and Completeness of  $\mathbf{TR}_1^{\text{bt}}$ .**

Suppose  $E$  is a bound-typed  $\lambda_1$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) It is possible to assign types of  $\lambda_1^{\text{type}}$  to all the variables of  $E$  in such a way that the free variables  $v_1, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \dots, \tau_n$ , and  $E$  itself has type  $X$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_1^{\text{bt}}$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

**Theorem 2.33: Subject Reduction Theorem for  $\mathbf{TR}_1$ .**

In  $\mathbf{TR}_1$ , if  $\Gamma \vdash M : \sigma$  and  $M \rightarrow_{\beta\eta} N$ , then  $\Gamma \vdash N : \sigma$ .

As in the previous chapter, many of these theorems have analogs for the bound-typed  $\lambda$ -calculus that we do not state separately, with the exception of the Soundness and Completeness Theorems 2.31 and 2.32.

One result that (perhaps surprisingly) does *not* generalize straightforwardly from  $\lambda_0$  to  $\lambda_1$  is the  $\eta$ -reduction postponement lemma (Lemma 1.24.) Consider the following reduction sequence in  $\lambda_1$

$$\langle \pi_1(\lambda x(x)), \pi_2(\lambda x(x)) \rangle(x) \rightarrow_{\eta} \lambda x(x)(x) \rightarrow_{\beta} x$$

The reduction from  $\langle \pi_1(\lambda x(x)), \pi_2(\lambda x(x)) \rangle(x)$  to  $x$  cannot be achieved by a sequence of  $\beta$ -reductions followed by a sequence of  $\eta$ -reductions. (Note that term  $\langle \pi_1(\lambda x(x)), \pi_2(\lambda x(x)) \rangle(x)$  is actually in  $\beta$ -normal form.) There is therefore no simple analog of Lemma 1.24 for  $\lambda_1$ . However, perhaps surprisingly it turns out that although  $\eta$ -reduction postponement fails for the *untyped* system  $\lambda_1$ , it holds for the *typed* system  $\lambda_1^{\text{type}}$ . (Note that the expression  $\langle \pi_1(\lambda x(x)), \pi_2(\lambda x(x)) \rangle$  cannot be properly typed, as  $\lambda x(x)$  has the type  $\tau \rightarrow \tau$  for some  $\tau$ , and  $\pi_1$  and  $\pi_2$  can only operate on terms of type  $\tau \times \sigma$ .)

In particular, we have

**Lemma 2.34:  $\eta$ -reduction postponement in  $\lambda_1^{\text{type}}$ .**

For all  $\lambda_1^{\text{type}}$ -terms  $X, Y, Z$ , if  $X \rightarrow_{\eta} Y \rightarrow_{\beta} Z$  in  $\lambda_1^{\text{type}}$ , then there is a  $\lambda_1^{\text{type}}$ -term  $Y'$  such that  $X \rightarrow_{\beta} Y' \rightarrow_{\eta} Z$  (where  $\rightarrow_{\eta}$  is a sequence of  $\eta$ -reductions.)

This is proved in the appendix. It follows from this that if a typed term can be reduced to  $\beta\eta$ -normal form, then there is a way of reducing it to  $\beta\eta$ -normal form that involves performing a sequence of  $\beta$ -reductions followed by a sequence of  $\eta$ -reductions.

**Exercises for Section 2.13**

1. For each of the following  $\lambda_1$ -terms, indicate if it can be typed or not. If it can be typed, write out the corresponding  $\lambda_1^{\text{type}}$  term. If it cannot, explain why not.

(i)  $\pi_1(x)(y)$

(ii)  $y(\pi_1(x))$

(iii)  $\pi_1(\lambda x(x))$

(iv)  $\langle x, y \rangle z$

(v)  $z \langle x, y \rangle$

(vi)  $\text{in}_1(x)(y)$

(vii)  $y(\text{in}_1(x))$

(viii)  $\text{case}(\langle x, y \rangle, [u]u, [v]v)$

(ix)  $\text{case}(x, [u]\langle u, u \rangle, [v]\lambda z(v))$

(x)  $\text{case}(x, [u]w(\langle u, u \rangle), [v]v)$

2. Prove Lemma 2.23

3. Verify that Corollaries 2.25 and 2.27 do in fact follow from their respective theorems.

4. Prove Theorem 2.30.

5. Prove Theorem 2.33.

**2.14 The Curry-Howard Correspondence: Version II**

If we take the rules of **TR**<sub>1</sub> and as before in each judgment  $x : \tau$  erase the  $\lambda$ -term  $x$  and keep only the type  $\tau$ , we get the following:

$$\begin{array}{c}
 \Gamma, \tau \vdash \tau \\
 \frac{\Gamma \vdash \tau \rightarrow \sigma \quad \Gamma \vdash \tau}{\Gamma \vdash \sigma} \\
 \frac{\Gamma, \tau \vdash \sigma}{\Gamma \vdash \tau \rightarrow \sigma} \\
 \frac{\Gamma \vdash \sigma \quad \Gamma \vdash \tau}{\Gamma \vdash \sigma \& \tau} \\
 \frac{\Gamma \vdash \sigma \& \tau}{\Gamma \vdash \sigma} \quad \frac{\Gamma \vdash \sigma \& \tau}{\Gamma \vdash \tau} \\
 \frac{\Gamma \vdash \sigma}{\Gamma \vdash \sigma \vee \tau} \quad \frac{\Gamma \vdash \tau}{\Gamma \vdash \sigma \vee \tau} \\
 \frac{\Gamma \vdash \sigma \vee \tau \quad \Gamma, \sigma \vdash \rho \quad \Gamma, \tau \vdash \rho}{\Gamma \vdash \rho}
 \end{array}$$

$$\frac{\Gamma \vdash \perp}{\Gamma \vdash \tau}$$

where  $\Gamma$  is an (unordered) list of formulae. Here,  $\times$  has been replaced by  $\&$ ,  $+$  has been replaced by  $\vee$ , and  $0$  has been replaced by  $\perp$ . What we have here are of course just the axioms and rules of inference for *full* intuitionistic propositional logic **NJ**.

For much of the remainder of this book, we will regard the symbols  $\&$  and  $\times$ , the symbols  $\vee$  and  $+$ , and the symbols  $0$  and  $\perp$  as interchangeable. For the sake of emphasis, in the context of a typing derivation or when discussing types we will (typically) use the symbols  $\times$ ,  $+$  and  $0$ , and in the context of a proof in logic we will (typically) use the symbols  $\&$ ,  $\vee$  and  $\perp$ . Thinking of these as just different ways of writing the same thing will allow us to avoid having to clutter notation with cumbersome syntactical ‘translation’ functions. Adopting this convention, starting with the rules of **TR**<sub>1</sub> and erasing all  $\lambda$ -terms and keeping only the types, we then *literally* get the rules of **NJ**.

The relationship between **TR**<sub>1</sub> and **NJ** is then similar to the relationship between **TR**<sub>0</sub> and **NJ**( $\rightarrow$ ). Every proof in **TR**<sub>1</sub> contains as its ‘skeleton’ a proof in **NJ**, and we will see that every proof in **NJ** is the ‘skeleton’ of some proof in **TR**<sub>0</sub>. In fact, we have the following more general form of the Curry-Howard Correspondence:

**Theorem 2.35: The Curry-Howard Correspondence for TR<sub>1</sub> and NJ, Part 1.**

If **TR**<sub>1</sub> proves

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

then **NJ** proves

$$A_1, \dots, A_n \vdash B$$

**Proof**

In any proof in **TR**<sub>1</sub>, if one takes each assertion of the form  $z : B$  and discards the term  $z$  leaving only the type  $B$ , the result is a proof of **NJ**.

**Theorem 2.36: The Curry-Howard Correspondence for TR<sub>1</sub> and NJ, Part 2.**

If **NJ** proves

$$A_1, \dots, A_n \vdash B$$

Then for any distinct variables  $x_1, \dots, x_n$ , there is some  $\lambda_1$ -term  $s$  such that **TR**<sub>1</sub> proves

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

**Proof**

The proof is a straightforward induction on the construction of proofs in **NJ**, and is completely analogous to the proof of Theorem 2.8. The base case, and the case in which the final step of the proof is an application of ( $\rightarrow$ E) or ( $\rightarrow$ I) are dealt with precisely as in the proof of Theorem 2.8.

It suffices to consider the case in which the final step of the **NJ** proof is an application of one of the rules ( $\&$ E), ( $\&$ I), ( $\vee$ I), ( $\vee$ E) or ( $\perp$ E).

For example, suppose that a proof of  $A_1, \dots, A_n \vdash B$  has as its final step an application of ( $\&$ I). This means that  $B$  has the form  $C\&D$  and the proof has the form:

$$\frac{\begin{array}{c} \vdots \\ A_1, \dots, A_n \vdash C \end{array} \quad \begin{array}{c} \vdots \\ A_1, \dots, A_n \vdash D \end{array}}{A_1, \dots, A_n \vdash C\&D} (\&I)$$

By inductive hypothesis we have proofs of both the following sequents:

$$x_1 : A_1, \dots, x_n : A_n \vdash s : C \quad \text{and} \quad x_1 : A_1, \dots, x_n : A_n \vdash t : D.$$

Combining these gives the proof

$$\frac{\begin{array}{c} \vdots \\ x_1 : A_1, \dots, x_n : A_n \vdash s : C \end{array} \quad \begin{array}{c} \vdots \\ x_1 : A_1, \dots, x_n : A_n \vdash t : D \end{array}}{x_1 : A_1, \dots, x_n : A_n \vdash \langle s, t \rangle : C \times D}$$

as desired.

In a similar way, proofs whose final steps are ( $\&$ E), ( $\vee$ I), ( $\vee$ E), or ( $\perp$ E)

$$\begin{array}{c} \frac{A_1, \dots, A_n \vdash C\&D}{A_1, \dots, A_n \vdash C} (\&E) \quad \frac{A_1, \dots, A_n \vdash C\&D}{A_1, \dots, A_n \vdash D} (\&E) \\ \frac{A_1, \dots, A_n \vdash C}{A_1, \dots, A_n \vdash C \vee D} (\vee I) \quad \frac{A_1, \dots, A_n \vdash D}{A_1, \dots, A_n \vdash C \vee D} (\vee I) \\ \frac{A_1, \dots, A_n \vdash C \vee D \quad A_1, \dots, A_n, C \vdash E \quad A_1, \dots, A_n, D \vdash E}{A_1, \dots, A_n \vdash E} (\vee E) \\ \frac{A_1, \dots, A_n \vdash \perp}{A_1, \dots, A_n \vdash C} (\perp E) \end{array}$$

are transformed to **TR**<sub>1</sub> typing derivations

$$\frac{\Gamma \vdash s : C \times D}{\Gamma \vdash \pi_1(s) : C} (\&E) \quad \frac{\Gamma \vdash s : C \times D}{\Gamma \vdash \pi_1(t) : D} (\&E)$$

$$\begin{array}{c}
\frac{\Gamma \vdash s : C}{\Gamma \vdash \mathbf{in}_1(s) : C + D} (\vee \text{I}) \quad \frac{\Gamma \vdash s : D}{\Gamma \vdash \mathbf{in}_2(s) : C + D} (\vee \text{I}) \\
\frac{\Gamma \vdash s : C + D \quad \Gamma, y : C \vdash s' : E \quad \Gamma, z : D \vdash s'' : E}{\Gamma \vdash \mathbf{case}(s, [y]s', [z]s'') : E} (\vee \text{E}) \\
\frac{\Gamma \vdash s : 0}{\Gamma \vdash \mathbf{abs}(s) : C} (\perp \text{E})
\end{array}$$

where  $\Gamma$  is  $x_1 : A_1, \dots, x_n : A_n$ .

These results show how we can move between proofs in intuitionistic logic **NJ** and typing derivations in **TR**<sub>1</sub>. Moving from a typing derivation in **TR**<sub>1</sub> to a proof in **NJ** is trivial; it suffices to simply delete  $\lambda$ -terms and keep types. Moving from a proof in **NJ** to a typing derivation in **TR**<sub>1</sub> is done as before, by simply assigning variables appropriately to the formulae (now viewed as types) in the leaves of the **NJ** proof, and then using the typing rules to generate the  $\lambda$ -terms progressively deeper in the proof.

To demonstrate this with a slightly more complex example, we have seen already that we can prove  $(\neg A \vee \neg B) \rightarrow \neg(A \& B)$  in **NJ** by combining the subproofs

$$\frac{A \& B, \neg A \vdash \neg A \quad \frac{A \& B, \neg A \vdash A \& B}{A \& B, \neg A \vdash A} (\& \text{I})}{A \& B, \neg A \vdash \perp} (\rightarrow \text{E})$$

and

$$\frac{A \& B, \neg B \vdash \neg B \quad \frac{A \& B, \neg B \vdash A \& B}{A \& B, \neg B \vdash B} (\& \text{I})}{A \& B, \neg B \vdash \perp} (\rightarrow \text{E})$$

to give

$$\frac{\frac{\neg A \vee \neg B \vdash \neg A \vee \neg B \quad \frac{\frac{\vdots \quad \vdots}{A \& B, \neg A \vdash \perp \quad A \& B, \neg B \vdash \perp} (\vee \text{E}')}{\neg A \vee \neg B, A \& B \vdash \perp} (\rightarrow \text{I})}{\neg A \vee \neg B \vdash \neg(A \& B)} (\rightarrow \text{I})}{\vdash (\neg A \vee \neg B) \rightarrow \neg(A \& B)} (\rightarrow \text{I})$$

The leaves of this larger proof consist of the formulae

$$\neg A \vee \neg B, A \& B, \neg A, \neg B$$

Let  $w, x, y, z$  be variables of these types, i.e.,

$$w : \neg A + \neg B, x : A \times B, y : \neg A, z : \neg B$$

The first subproof is transformed into the typing derivation



$$\frac{x : A \times B, y : \neg A \vdash y : \neg A \quad \frac{x : A \times B, y : \neg A \vdash x : A \times B}{x : A \times B, y : \neg A \vdash \pi_1(x) : A}}{x : A \times B, y : \neg A \vdash y(\pi_1(x)) : \perp}$$

The second subproof is likewise transformed into the typing derivation

$$\frac{x : A \times B, z : \neg B \vdash z : \neg B \quad \frac{x : A \times B, z : \neg B \vdash x : A \times B}{x : A \times B, z : \neg B \vdash \pi_2(x) : B}}{x : A \times B, z : \neg B \vdash z(\pi_2(x)) : \perp}$$

putting these together yields:

$$\frac{w : \neg A + \neg B \vdash w : \neg A + \neg B \quad \begin{array}{c} \vdots \\ x : A \times B, y : \neg A \vdash y(\pi_1(x)) : \perp \\ \vdots \end{array} \quad \begin{array}{c} \vdots \\ x : A \times B, z : \neg B \vdash z(\pi_2(x)) : \perp \\ \vdots \end{array}}{w : \neg A + \neg B, x : A \times B \vdash \mathbf{case}(w, [y]y(\pi_1(x)), [z]z(\pi_2(x))) : \perp}}{\frac{w : \neg A + \neg B \vdash \lambda x \mathbf{case}(w, [y]y(\pi_1(x)), [z]z(\pi_2(x))) : \neg(A \times B)}{\vdash \lambda w \lambda x \mathbf{case}(w, [y]y(\pi_1(x)), [z]z(\pi_2(x))) : (\neg A + \neg B) \rightarrow \neg(A \times B)}}$$

where for the sake of compactness typing declarations have been omitted from contexts when not needed. This of course is a typing derivation in  $\mathbf{TR}_1$ .

As before, we can use this version of the Curry-Howard correspondence to solve technical problems about  $\mathbf{TR}_1$  by using intuitionistic logic. For example, for arbitrary types  $A$  and  $B$  the type  $(A + B) \rightarrow A$  is *not* inhabited, as  $(A \vee B) \rightarrow A$  is not a theorem of intuitionistic (or classical) logic. In this way, understanding the structure of  $\mathbf{NJ}$  gives us a tool for understanding the structure of  $\mathbf{TR}_1$ .

Conversely, understanding the structure of  $\mathbf{TR}_1$  can help us to understand the structure of  $\mathbf{NJ}$ . We begin with a more superficial and familiar type of example. Is  $A \rightarrow (B \rightarrow (A \& B))$  intuitionistically provable? This question can be answered easily enough by examining the rules of  $\mathbf{NJ}$ . But let us consider this question through the lens of the Curry-Howard correspondence. By the Curry-Howard correspondence, our question is equivalent to whether we can find a term in our expanded  $\lambda$ -calculus that can be assigned the type  $A \rightarrow (B \rightarrow (A \times B))$ . In fact, it is not hard to see that  $\lambda x \lambda y (\langle x, y \rangle)$  is such a term. If  $x$  has type  $A$  and  $y$  has type  $B$ , then this term has type  $A \rightarrow (B \rightarrow (A \times B))$ , as the following typing derivation shows:

$$\frac{\frac{x : A, y : B \vdash x : A \quad x : A, y : B \vdash y : B}{x : A, y : B \vdash \langle x, y \rangle : A \times B}}{x : A \vdash \lambda y (\langle x, y \rangle) : B \rightarrow (A \times B)}{\vdash \lambda x \lambda y (\langle x, y \rangle) : A \rightarrow (B \rightarrow (A \times B))}$$

Deleting  $\lambda$ -terms then yields the following proof of  $A \rightarrow (B \rightarrow (A \& B))$  in  $\mathbf{NJ}$

$$\frac{\frac{A, B \vdash A \quad A, B \vdash B}{A, B \vdash A \& B}}{A \vdash B \rightarrow (A \& B)}{\vdash A \rightarrow (B \rightarrow (A \& B))}$$

We will see more substantive applications of  $\mathbf{TR}_1$  to our understanding of  $\mathbf{NJ}$  in the next sections when we return to the topic of proof normalization.

### Exercises for Section 2.14

1. The following sentences  $X$  are all theorems of  $\mathbf{NJ}$ . For each of them, first construct a proof in  $\mathbf{NJ}$ . (You have already constructed proofs for some of these in earlier exercises.) Then transform each of these proofs of  $X$  into a typing derivation in  $\mathbf{TR}_1$  of  $\vdash M : X$  for some  $\lambda_1$ -term  $M$ .

- (a)  $p \rightarrow (p \vee q)$
- (b)  $\neg(p \vee q) \rightarrow \neg p$
- (c)  $(p \& \neg p) \rightarrow q$
- (d)  $((p \rightarrow r) \& (q \rightarrow r)) \rightarrow ((p \vee q) \rightarrow r)$
- (e)  $\neg\neg(p \vee \neg p)$

2. Are all the  $\lambda_1$ -terms  $M$  obtained in problem 1.(a)-(e) in  $\beta$ -normal form? If not, in each case find a different  $\mathbf{NJ}$  proof for which the relevant  $\lambda_1$ -term  $M$  is in  $\beta$ -normal form.

## 2.15 Proof Normalization Revisited

Earlier, we examined  $\mathbf{NJ}(\rightarrow)$  proofs in which an application of the ( $\rightarrow$ I) rule was immediately followed by an application of the ( $\rightarrow$ E) rule, with the conclusion of the ( $\rightarrow$ I) rule used as the major premise:

$$\frac{\frac{\Gamma, A \vdash B}{\Gamma \vdash A \rightarrow B} (\rightarrow\text{I}) \quad \Gamma \vdash A}{\Gamma \vdash B} (\rightarrow\text{E})$$

The typing derivations corresponding to such proofs look like:

$$\frac{\frac{\bar{\Gamma}, u : A \vdash e : B}{\bar{\Gamma} \vdash \lambda u(e) : A \rightarrow B} \quad \bar{\Gamma} \vdash f : A}{\bar{\Gamma} \vdash \lambda u(e)(f) : B}$$

The term  $\lambda u(f)(g)$  is of course  $\beta$ -reducible. We concluded that typing derivations of  $\beta$ -normal form terms corresponded to  $\mathbf{NJ}(\rightarrow)$  proofs without the pattern of reasoning in question. This was formally stated as Theorem 2.13.

This basic idea generalizes to our other connectives. Consider a proof in which an application of the ( $\&$ I) rule is immediately followed by an application of the ( $\&$ E) rule:

$$\frac{\frac{\Gamma \vdash A \quad \Gamma \vdash B}{\Gamma \vdash A \& B} (\&\text{I})}{\Gamma \vdash A} (\&\text{E})$$

A typing derivation corresponding to this proof will look like:

$$\frac{\frac{\bar{\Gamma} \vdash e : A \quad \bar{\Gamma} \vdash f : B}{\bar{\Gamma} \vdash \langle e, f \rangle : A \times B}}{\bar{\Gamma} \vdash \pi_1(\langle e, f \rangle) : A}$$

Note that in this case too the final term  $\pi_1(\langle e, f \rangle)$  is not in  $\beta$ -normal form.

Finally, consider a proof in which an application of the ( $\vee$ I) rule is immediately followed by an application of the ( $\vee$ E) rule, with the conclusion of the ( $\vee$ I) rule used as the major premise. Here there are two possibilities:

$$\frac{\frac{\Gamma \vdash A}{\Gamma \vdash A \vee B} (\vee\text{I}) \quad \Gamma, A \vdash C \quad \Gamma, B \vdash C}{\Gamma \vdash C} (\vee\text{E})$$

and

$$\frac{\frac{\Gamma \vdash B}{\Gamma \vdash A \vee B} (\vee\text{I}) \quad \Gamma, A \vdash C \quad \Gamma, B \vdash C}{\Gamma \vdash C} (\vee\text{E})$$

The typing derivations corresponding to these will look like:

$$\frac{\frac{\bar{\Gamma} \vdash e : A}{\bar{\Gamma} \vdash \mathbf{in}_1(E) : A + B} \quad \bar{\Gamma}, u : A \vdash f : C \quad \bar{\Gamma}, v : B \vdash g : C}{\bar{\Gamma} \vdash \mathbf{case}(\mathbf{in}_1(e), [u]f, [v]g) : C}}$$

and

$$\frac{\frac{\bar{\Gamma} \vdash e : B}{\bar{\Gamma} \vdash \mathbf{in}_2(E) : A + B} \quad \bar{\Gamma}, u : A \vdash f : C \quad \bar{\Gamma}, v : B \vdash g : C}{\bar{\Gamma} \vdash \mathbf{case}(\mathbf{in}_2(e), [u]f, [v]g) : C}}$$

Again, the final term  $\mathbf{case}(\mathbf{in}_1(e), [u]f, [v]g)$  or  $\mathbf{case}(\mathbf{in}_2(e), [u]f, [v]g)$  is not in  $\beta$ -normal form.

With this in mind, we consider the following definition

**Definition 2.37: Normal Proof of NJ.**

A proof in **NJ** is *normal* just in case

- (i) it contains no occurrence of a ( $\rightarrow$ I) rule introducing a conditional  $X \rightarrow Y$  which is immediately followed by an application of the ( $\rightarrow$ E) rule with this conditional  $X \rightarrow Y$  as its major premise.
- (ii) it contains no occurrence of a ( $\&$ I) rule immediately followed by an application of the ( $\&$ E) rule.
- (iii) it contains no occurrence of a ( $\vee$ I) rule introducing a disjunction  $X \vee Y$  which is immediately followed by an application of the ( $\vee$ E) rule with this disjunction  $X \vee Y$  as its

major premise.

As before we have

**Lemma 2.38: The Subterm Lemma for  $\mathbf{TR}_1$**

In any typing derivation of  $\mathbf{TR}_1$  with conclusion of the form  $\Gamma \vdash s : A$ , if  $s'$  is a subterm of  $s$  then the typing derivation contains a sequent of the form  $\Gamma' \vdash s' : B$ , where  $\Gamma' \subseteq \Gamma$ .

and

**Theorem 2.39**

Given a typing derivation  $T$  in  $\mathbf{TR}_1$  whose conclusion is a sequent of the form  $\Delta \vdash g : C$ , let  $P$  be the  $\mathbf{NJ}$  proof obtained from  $T$  by deleting all terms and keeping types. Then  $P$  is a normal proof iff  $g$  is in  $\beta$ -normal form.

As before, we may conclude

**Theorem 2.40: Proof Normalization Theorem for  $\mathbf{NJ}$**

If  $A_1, \dots, A_n \vdash B$  has a proof in  $\mathbf{NJ}$ , then it has a normal proof in  $\mathbf{NJ}$ .

**Proof**

Suppose  $P_1$  is a non-normal proof of  $A_1, \dots, A_n \vdash B$ . By the Curry-Howard Correspondence (in particular, Theorem 2.36), this may be transformed into a typing derivation with conclusion

$$x_1 : A_1, \dots, x_n : A_n \vdash s : B$$

for some  $\lambda$ -term  $s$ . By Theorem 2.39,  $s$  is not in  $\beta$ -normal form. Let  $s'$  be the  $\beta$ -normal form of  $s$ . By Theorem 2.33 (the Subject Reduction Theorem) there is a typing derivation of

$$x_1 : A_1, \dots, x_n : A_n \vdash s' : B$$

in  $\mathbf{TR}_1$ . Because  $s'$  is in  $\beta$ -normal form, it follows from Theorem 2.39 that the  $\mathbf{NJ}$  proof obtained by taking this typing derivation and removing all terms and keeping types is a normal proof of  $A_1, \dots, A_n \vdash B$ .

As before, this gives us an algorithm for proof normalization. Consider for example the following non-normal proof of  $A \rightarrow (B \rightarrow A)$  in  $\mathbf{NJ}$

$$\frac{\frac{A, B \vdash A}{A \vdash B \rightarrow A} (\rightarrow \text{I}) \quad A \vdash A \text{ } (\& \text{I})}{\frac{A \vdash (B \rightarrow A) \& A}{A \vdash B \rightarrow A} (\& \text{E})} (\rightarrow \text{I})$$

This is non-normal because it involves the (& I) rule immediately followed by the (& E) rule. This has the corresponding typing derivation

$$\frac{\frac{x : A, y : B \vdash x : A}{x : A \vdash \lambda y(x) : B \rightarrow A} \quad x : A \vdash x : A}{\frac{x : A \vdash \langle \lambda y(x), x \rangle : (B \rightarrow A) \times A}{x : A \vdash \pi_1 \langle \lambda y(x), x \rangle : B \rightarrow A}}{\vdash \lambda x(\pi_1 \langle \lambda y(x), x \rangle) : A \rightarrow (B \rightarrow A)}$$

The term  $\lambda x(\pi_1 \langle \lambda y(x), x \rangle)$  in the conclusion is not in  $\beta$ -normal form. It  $\beta$ -reduces to  $\lambda x \lambda y(x)$ , which is in  $\beta$ -normal form. The term  $\lambda x \lambda y(x)$  may be given the type  $A \rightarrow (B \rightarrow A)$  as follows

$$\frac{x : A, y : B \vdash x : A}{\frac{x : A \vdash \lambda y(x) : B \rightarrow A}{\vdash \lambda x \lambda y(x) : A \rightarrow (B \rightarrow A)}}$$

which corresponds to the **NJ** proof

$$\frac{A, B \vdash A}{\frac{A \vdash B \rightarrow A}{\vdash A \rightarrow (B \rightarrow A)} (\rightarrow \text{I})} (\rightarrow \text{I})$$

This proof is indeed normal.

Let us turn to an examination of the structure of normal proofs. Recall the definition of a *major branch*, which requires no changes

**Definition 2.41: Major Branch**

A *major branch* of a proof is a branch of the proof with the property that if the conclusion of an inference lies on the branch and that inference has a major premise, then that major premise also lies on in the branch.

As before, we easily see that each proof has at least one major branch. We also have the following:

**Lemma 2.42: The Major Branch Lemma.**

In any normal proof in **NJ**, a major branch will consist of a (possibly empty) sequence of elimination rules (i.e.,  $(\rightarrow E)$ ,  $(\& E)$ ,  $(\vee E)$  and  $(\perp E)$  rules) followed by a (possibly empty) sequence of introduction rules (i.e.,  $(\rightarrow I)$ ,  $(\& I)$  and  $(\vee I)$  rules.)

**Proof**

We first argue that along a major branch in a normal proof, a  $(\rightarrow I)$  rule cannot be immediately followed by *any* elimination rule. To see why, note that because the conclusion of a  $(\rightarrow I)$  rule has the form  $\Delta \vdash \phi \rightarrow \psi$ , it cannot be immediately followed by a  $(\& E)$  or a  $(\perp E)$  rule. Because the only premise of a  $(\vee E)$  rule that appears on a major branch is the major premise, which always has the form  $\Gamma \vdash \phi \vee \phi$ , it cannot be followed by a  $(\vee E)$  rule either. If it is followed by an elimination rule, it must therefore be followed by a  $(\rightarrow E)$  rule. But this happening would contradict the fact that the proof is normal. So a  $(\rightarrow I)$  rule cannot be followed by *any* elimination rule.

In a precisely similar way, it may be argued that along a major branch, a  $(\& I)$  rule cannot be followed by any elimination rule, and that a  $(\vee I)$  rule cannot be followed by any elimination rule. It thus follows that along a major branch of a normal proof, an introduction rule can only be followed by another introduction rule. It follows that any major branch of a normal must consist of a (possibly empty) sequence of elimination rules followed by a (possibly empty) sequence of introduction rules.

As before, the Major Branch Lemma gives a wealth of information about what is provable in **NJ**. For example:

**Theorem 2.43: Consistency of NJ**

No atomic sentence  $X$  is provable in **NJ**.

Recall that the only atomic sentences in our language are the propositional variables  $X, Y, Z, \dots$  and the sentence  $\perp$ . The proof is essentially identical to the consistency proof given for **NJ**( $\rightarrow$ ) (i.e., Theorem 2.17).

**Proof**

Suppose **NJ** proves  $S$  where  $S$  is either a propositional variable or the symbol  $\perp$ . We then know from Theorem 2.40 that there is a normal proof  $P$  whose conclusion is  $\vdash S$ . The last step of  $P$  cannot be an introduction rule, as each introduction rule introduces a non-atomic

sentence into the conclusion. The last step of the proof is therefore an elimination rule. Using the Major Branch Lemma (Lemma 2.42), it then follows that any major branch of  $P$  must consist solely of elimination rules. But along a major branch, the elimination rules all leave the left hand side of sequents unchanged. Because the left hand side of every leaf is non-empty, it follows that the conclusion of  $P$  must have a non-empty left hand side, which is a contradiction.

The following result is also of interest

**Theorem 2.44**

The only atomic sentence  $X$  for which  $\neg X$  provable in **NJ** is  $\perp$ .

**Proof**

Suppose **NJ** proves  $\neg X$  where  $X$  is an atomic sentence distinct from  $\perp$ . Consider a normal proof of  $\vdash \neg X$ . A major branch of this proof could not consist entirely of elimination rules, because the left hand side of  $\vdash \neg X$  is empty. Therefore the proof must end in an introduction rule. Because  $\neg X$  is a conditional, the only possible introduction rule the proof could end in is ( $\rightarrow$  I). The proof therefore has the form:

$$\frac{\vdots}{X \vdash \perp} (\rightarrow \text{I})$$

The sequent  $X \vdash \perp$  cannot be the conclusion of an introduction rule, and so must be the conclusion of an elimination rule. Thus, any major branch must consist entirely of elimination rules, followed by the ( $\rightarrow$  I) inference at the end of the proof. Because elimination rules do not change the left hand side of a sequent along a major branch, the proof must have the form

$$\frac{X \vdash X}{\vdots} (\rightarrow \text{I})$$

For atomic sentences  $X$ , the only case in which an elimination rule can be performed on the leaf  $X \vdash X$  on a major branch is the case in which  $X$  is  $\perp$  (in which case a ( $\perp$  E) can be performed). Because we are assuming that  $X$  is distinct from  $\perp$ , it follows that no elimination rule can be performed on the leaf  $X \vdash X$ , and thus it is impossible to derive  $X \vdash \perp$  solely using elimination rules. On the other hand, if  $X$  is  $\perp$ , then

$$\frac{\perp \vdash \perp}{\vdash \neg \perp} (\rightarrow \text{I})$$

is a derivation of  $\neg X$  in **NJ**.

From this, we easily have

**Theorem 2.45: Underivability of Double Negation Elimination**

If  $X$  is atomic sentence distinct from  $\perp$ , then  $\neg\neg X \rightarrow X$  is not provable in **NJ**.

**Proof**

Suppose to the contrary that  $\neg\neg X \rightarrow X$  is provable in **NJ** for some atomic sentence  $X$  distinct from  $\perp$ , and consider a normal proof of  $\vdash \neg\neg X \rightarrow X$ . The last step in this proof cannot be an elimination rule, as this would mean that any major branch of this proof would consist entirely of elimination rules, and thus the concluding sequent would have non-empty left hand side. So the final rule must be an introduction rule. Given that  $\neg\neg X \rightarrow X$  is neither a conjunction nor a disjunction, the only possibility is that the last rule is ( $\rightarrow$  I):

$$\frac{\neg\neg X \vdash X}{\vdash \neg\neg X \rightarrow X} (\rightarrow \text{I})$$

Because the right hand side of the sequent  $\neg\neg X \vdash X$  is atomic, it cannot be the result of an introduction rule, and must be the result of an elimination rule. Any major branch of the proof must therefore consist entirely of elimination rules followed by a single ( $\rightarrow$  I). Because along a major branch the left hand side of a sequent remains unchanged when performing elimination rules, any major branch must therefore have the form

$$\begin{array}{c} \neg\neg X \vdash \neg\neg X \\ \vdots \\ \frac{\neg\neg X \vdash X}{\vdash \neg\neg X \rightarrow X} (\rightarrow \text{I}) \end{array}$$

The only elimination rule that can be applied to the leaf  $\neg\neg X \vdash \neg\neg X$  along a major branch is ( $\rightarrow$  E), which means that  $\neg\neg X$  (i.e.,  $(X \rightarrow \perp) \rightarrow \perp$ ) must be the major premise, and the proof has the form

$$\begin{array}{c} \vdots \\ \frac{\frac{\neg\neg X \vdash \neg\neg X}{\neg\neg X \vdash \perp} \quad \neg\neg X \vdash \neg\neg X}{\vdash \neg\neg X \rightarrow \perp} (\rightarrow \text{E}) \\ \vdots \\ \frac{\neg\neg X \vdash X}{\vdash \neg\neg X \rightarrow X} (\rightarrow \text{I}) \end{array}$$



Thus we have that  $\neg\neg X \vdash \perp$  is derivable in **NJ**. Using  $(\rightarrow I)$  we then have that  $\vdash \neg\neg\neg X$  is derivable in **NJ**, and using the fact that  $\vdash \neg\neg\neg X \rightarrow \neg X$  is derivable in **NJ** and  $(\rightarrow E)$ , we finally have that  $\vdash \neg X$  is derivable in **NJ**. This however contradicts Theorem 2.44.

Note that the sequent  $\neg\neg \perp \rightarrow \perp$  is indeed provable in **NJ**. The proof is

$$\frac{\frac{\frac{\neg\neg \perp \vdash \neg\neg \perp}{\neg\neg \perp \vdash \perp} (\rightarrow I)}{\neg\neg \perp \vdash \neg\neg \perp} (\rightarrow E)}{\vdash \neg\neg \perp \rightarrow \perp} (\rightarrow I)$$

Some further applications of the Major Branch Lemma will be given in the exercises. One of the most interesting things about our proofs of Theorems 2.43, 2.44 and 2.45 is that they show the unprovability of certain sequents on purely syntactic grounds. No semantic considerations of any sort are invoked. That is to say, we do not construct ‘models’ in which certain sentences are false, thereby somehow showing that these sentences cannot be derived. Instead, our proofs merely consider the syntactic form of the sentences in question to argue for their unprovability.

### Exercises for Section 2.15

1. Each of the following proofs in **NJ** is *not* normal. Use the techniques of this section to transform them into normal proofs.

(a)

$$\frac{\frac{\frac{A, B, C, C \rightarrow A \vdash C}{B, C, C \rightarrow A \vdash C \rightarrow A} (\rightarrow I)}{B, C, C \rightarrow A \vdash A} (\rightarrow E)}{\frac{B, C, C \rightarrow A \vdash A \& B}{B, C, C \rightarrow A \vdash A} (\&E)} \quad \frac{B, C, C \rightarrow A \vdash B}{B, C, C \rightarrow A \vdash A \& B} (\&I)$$

(b)

$$\frac{\frac{\frac{A \& B \vdash A \& B}{A \& B \vdash A} (\&E)}{A \& B \vdash A \vee B} (\vee I)}{A \& B \vdash B} \quad \frac{A \& B, A \vdash A \& B}{A \& B, A \vdash B} (\&E)}{\frac{A \& B, B \vdash B}{A \& B, B \vdash A \& B} (\&I)}{A \& B, B \vdash B} (\vee E)$$

2. Using the Major Branch Lemma, show that the following sequents are not derivable in **NJ**, where  $X$  and  $Y$  are *distinct* atomic sentences, with neither equal to  $\perp$ .

(i)  $X \vdash Y$

(ii)  $\vdash X \vee \neg X$

(iii)  $\neg(X \& Y) \vdash X$ , and likewise  $\neg(X \& Y) \vdash Y$ . (Hint: consider the normal proof with the fewest number of sequents.)

(iv)  $\neg(X \& Y), X \vdash Y$

(v)  $\neg(X \& Y), X \vdash \perp$ , and likewise  $\neg(X \& Y), Y \vdash \perp$

(vi)  $\vdash \neg(X \& Y) \rightarrow (\neg X \vee \neg Y)$

3. By using the Major Branch Lemma, show that Pierce's Rule is not derivable in **NJ**. (You should be able to adapt the proof that it is not provable in **NJ**( $\rightarrow$ ).)

4. Show that if  $X$  is a sentence constructed from atomic sentences (including  $\perp$ ) using only the connective  $\rightarrow$ , then if **NJ** proves  $\vdash X$ , then **NJ**( $\rightarrow$ ) proves  $\vdash X$ .

## 2.16 The BHK Interpretation Revisited

We finally return to topic of the **BHK** conditions for intuitionistic logic. Armed with everything we have presented so far, we will see that there is a sense in which these laws can be vindicated in a concrete and natural way.

Recall the four **BHK** conditions:

**BHK1** : To prove  $A \& B$  is to prove  $A$  and to prove  $B$ .

**BHK2** : To prove  $A \vee B$  is to prove  $A$  or to prove  $B$ .

**BHK3** : To prove  $A \rightarrow B$  is to construct a way of transforming a proof of  $A$  into a proof of  $B$ .

**BHK4** : There is no proof of  $\perp$ .

There are two main obstacles in making rigorous sense of these claims. To see the first, let us consider **BHK1**. It is *not* true that in intuitionistic logic, any proof of  $A \& B$  literally consists of a proof of  $A$  and a proof of  $B$ , conjoined into a proof of  $A \& B$  by the ( $\&I$ ) rule. This is because we can introduce unnecessary detours into a proof. For example, a proof of  $A$  and a proof of  $B$  can lead to a proof of  $A \& B$  as follows:

$$\frac{\frac{\frac{\vdash A \quad \vdash B}{\vdash A \& B} (\&I) \quad \vdash A}{\vdash (A \& B) \& A} (\&I)}{\vdash A \& B} (\&E)$$

The problem with this proof of course is that it is not normal. A natural idea then is to think of the **BHK** conditions as claims about *normal* proofs.

A second obstacle in making sense of the **BHK** conditions revolves around **BHK3**. This claim tells us that proofs of conditionals are *functions* of a certain sort. But **NJ** proofs are not functions in any literal sense – they are trees of sequents. We turn to the question of how to make sense of **BHK3** shortly.

We begin with the following theorem

**Theorem 2.46: BHK Theorem**

The following are all true

(i) If a normal proof in **NJ** has a conclusion of the form  $\vdash X \& Y$ , then it must have the form

$$\frac{\begin{array}{c} \vdots \\ \vdash X \end{array} \quad \begin{array}{c} \vdots \\ \vdash Y \end{array}}{\vdash X \& Y} \text{ (& I)}$$

where the proofs of the sequents  $\vdash X$  and  $\vdash Y$  are also normal.

(ii) If a normal proof in **NJ** has a conclusion of the form  $\vdash X \vee Y$ , then it must have one of the forms:

$$\frac{\begin{array}{c} \vdots \\ \vdash X \end{array}}{\vdash X \vee Y} (\vee \text{ I}) \quad \text{or} \quad \frac{\begin{array}{c} \vdots \\ \vdash Y \end{array}}{\vdash X \vee Y} (\vee \text{ I})$$

where the proof of the sequent  $\vdash X$  or  $\vdash Y$  is also normal.

(iii) If a normal proof in **NJ** has a conclusion of the form  $\vdash X \rightarrow Y$ , then it must have the form

$$\frac{\begin{array}{c} \vdots \\ X \vdash Y \end{array}}{\vdash X \rightarrow Y} (\rightarrow \text{ I})$$

where the proof of the sequent  $X \vdash Y$  is also normal.

(iv) No normal proof in **NJ** has the form

$$\begin{array}{c} \vdots \\ \vdash \perp \end{array}$$

**Proof**

For (i), suppose a normal proof in **NJ** has a conclusion of the form  $\vdash X \& Y$ , and consider a major branch of this proof. This major branch cannot consist entirely of elimination rules, as the conclusion  $\vdash X \& Y$  has an empty left hand side. Therefore the final rule of inference in this branch must be an introduction rule. Because  $X \& Y$  has the form of a conjunction, the only introduction rule which could prove the sequent  $\vdash X \& Y$  is the (& I) rule. The proof must therefore have the form

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \quad \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array}}{\vdash X \quad \vdash Y} (\& I)$$

Because subproofs of normal proofs are also normal proofs, the proofs of the sequents  $\vdash X$  and  $\vdash Y$  are also normal.

The proofs of (ii) and (iii) involve essentially the same form of argument. That (iv) is true follows from Theorem 2.43.

Theorem 2.46 shows that we may view **BHK1**, **BHK2** and **BHK4** as true claims about *normal* proofs.

Let us turn to condition **BHK3**. Is there a sense in which a normal proof of  $A \rightarrow B$  is a function that takes us from a normal proof of  $A$  to a normal proof of  $B$ ? While a normal proof of  $A \rightarrow B$  is not *literally* a function from normal proofs of  $A$  to normal proofs of  $B$ , it is easily seen to encode such a function. For given any proofs of  $A \rightarrow B$  and  $A$  in **NJ**, we can produce a proof of  $B$  in **NJ** by invoking ( $\rightarrow$  E):

$$\frac{\vdash A \rightarrow B \quad \vdash A}{\vdash B} (\rightarrow E)$$

This proof need not be normal, but normalizing it in the way described in the previous sections yields a normal proof. Thus, given any proof of  $A \rightarrow B$ , we have an algorithm that takes a normal proof of  $A$  to a normal proof of  $B$ . In this sense, **BHK3** is also vindicated.

Thus we see that all of **BHK1-4** may be interpreted as true claims about normal proofs of **NJ**. In this sense, the proof system **NJ** does indeed capture the ideas behind the **BHK** interpretation, and we are therefore right to think of **NJ** as a formalization of ‘constructive’ or ‘intuitionistic’ reasoning. One advantage (among many) of the formalism of type theory and the Curry-Howard correspondence is the way in which it allows us to see all of this particularly clearly.

### Exercises for Section 2.16

1. Show the following.

- (i) It is *not* true in general that if **NJ** proves a sequent of the form  $\Delta \vdash X \vee Y$ , then it must prove either  $\Delta \vdash X$  or  $\Delta \vdash Y$ .
- (ii) If each sentence in  $\Delta$  has the form  $\neg Z$ , then if **NJ** proves  $\Delta \vdash X \vee Y$ , it must also prove either  $\Delta \vdash X$  or  $\Delta \vdash Y$ .

## 2.17 Other Resources

For a good discussion of propositional intuitionistic logic and the Curry-Howard Correspondence, see chapters 2 and 4 of Sorenson and Urzyczyn’s *Lectures on the Curry-Howard Isomorphism* [10] and chapter 6 of Selinger’s *Lecture Notes on the Lambda Calculus* [9]. Another good reference focusing

mainly on the connection with intuitionistic logic is Mints' *A Short Introduction to Intuitionistic Logic* [7]. For a discussion of the themes of this chapter more from the point of view of computer science, chapters 8 through 14 of Pierce's *Types and Programming Languages* [8] is a good introduction.

## 2.18 Appendix

In this section, we provide proofs for several results in the main text.

### Lemma 2.34: $\eta$ -reduction postponement in $\lambda_1$ .

For all  $\lambda_1^{\text{type}}$ -terms  $X, Y, Z$ , if  $X \rightarrow_\eta Y \rightarrow_\beta Z$  in  $\lambda_1^{\text{type}}$ , then there is a  $\lambda_1^{\text{type}}$ -term  $Y'$  such that  $X \rightarrow_\beta Y' \rightarrow_\eta Z$  (where  $\rightarrow_\eta$  is a sequence of  $\eta$ -reductions.)

### Proof of Lemma 2.34.

As very little of the argument explicitly uses the fact that we are working in the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$ , we will omit typing superscripts in what follows for ease of readability, instead explicitly drawing attention to facts about typing when needed.

Suppose that we have  $\lambda_1^{\text{type}}$ -terms  $X, Y, Z$  with  $X \rightarrow_\eta Y \rightarrow_\beta Z$ . Then  $X$  can be written in the form  $Et'F$  and  $Y$  in the form  $EtF$  where  $E$  and  $F$  are sequences of symbols, and  $t'$  is an  $\eta$ -redex – i.e., a term with one of the forms  $\lambda x(t(x))$  (with  $x$  not free in  $t$ ),  $\langle \pi_1(t), \pi_2(t) \rangle$ , or  $\text{case}(t, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v))$ . Thus we write

$$Et'F \rightarrow_\eta EtF \rightarrow_\beta Z$$

where  $Z$  is a  $\lambda$ -term.

The  $\beta$ -reduction  $EtF \rightarrow_\beta Z$  occurs by reducing some  $\beta$ -redex  $s$  which occurs a subexpression of  $EtF$ . We then have either that (i)  $s$  is a subexpression of  $t$  (where this includes the possibility that  $s = t$ ), (ii)  $t$  is a proper subexpression of  $s$ , or (iii)  $s$  and  $t$  are disjoint.

Consider first case (i). Assume that when the  $\beta$ -reduction in question is performed on the subexpression  $s$  of  $t$ , we get  $t \rightarrow_\beta t''$ . We can then represent our reduction sequence as

$$Et'F \rightarrow_\eta EtF \rightarrow_\beta Et''F$$

Because  $t \rightarrow_\beta t''$ , we may then reverse the order of the  $\beta$ - and  $\eta$ -reductions for each of our 3 possible types of  $\eta$ -reduction as follows

$$Et'F = E \begin{array}{c} \lambda x(t(x)) \\ \langle \pi_1(t), \pi_2(t) \rangle \\ \text{case}(t, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) \end{array} F \rightarrow_\beta E \begin{array}{c} \lambda x(t''(x)) \\ \langle \pi_1(t''), \pi_2(t'') \rangle \\ \text{case}(t'', [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v)) \end{array} F \rightarrow_\eta Et''F$$

In the case in which  $t'$  is  $\lambda x(t(x))$ , we must use the fact that  $x$  is not free in  $t''$ . This is true because  $x$  is not free in  $t$  and no  $\beta$ -reduction introduces a new free variable.

Consider next case (iii). If  $s$  and  $t$  are disjoint, then  $s$  is either a subexpression of  $E$  or a subexpression of  $F$ . Suppose without loss of generality that  $s$  is a subexpression of  $E$ , so that  $E$  can be written as  $GsH$ . Then our reduction sequence has the form

$$GsHt'F \rightarrow_{\eta} GsHtF \rightarrow_{\beta} Gs'HtF$$

where  $s \rightarrow_{\beta} s'$  and  $t' \rightarrow_{\eta} t$ . In this case, the order of the  $\beta$ - and  $\eta$ -reductions may easily be reversed as follows

$$GsHt'F \rightarrow_{\beta} Gs'Ht'F \rightarrow_{\eta} Gs'HtF$$

Case (ii) ( $t$  is a proper subexpression of  $s$ ) requires the most care. In this case we have

$$Et'F \rightarrow_{\eta} EtF \rightarrow_{\beta} Z$$

where the  $\beta$ -redex  $s$  being reduced in  $EtF \rightarrow_{\beta} Z$  has  $t$  as a proper subexpression. We consider each of the possible form of  $\beta$ -redex in turn.

First, suppose  $s$  has the form  $\lambda z(M)(N)$ . Because  $t$  is a proper subexpression of  $s$ , we have that (a)  $t$  is  $\lambda z(M)$ , (b)  $t$  is a subexpression of  $M$ , or (c)  $t$  is a subexpression of  $N$ .

In case (a) we use some facts about typing. (This is one of the three points in the proof that relies on the fact that we are working in the *typed* setting.) If  $t$  is  $\lambda z(M)$ , then the type of  $t$  has the form  $\alpha \rightarrow \beta$ . The only  $\eta$ -reduction that returns an object of such a type is the reduction rule  $\lambda x(E(x)) \rightarrow_{\eta} E$ . Thus,  $t'$  has the form  $\lambda x(\lambda z(M)(x))$  and the reduction sequence has the form

$$E\lambda x(\lambda z(M)(x))F \rightarrow_{\eta} E\lambda z(M)F \rightarrow_{\beta} Z$$

where  $F$  has the form  $(N)F'$ , i.e.,

$$E\lambda x(\lambda z(M)(x))(N)F' \rightarrow_{\eta} E\lambda z(M)(N)F' \rightarrow_{\beta} EM[N/z]F'$$

where  $x$  does not occur free in  $\lambda z(M)$ . Assuming as we may that  $x$  and  $z$  are distinct variables, this means that  $x$  does not occur free in  $M$ . But then we also have the reduction sequence

$$E\lambda x(\lambda z(M)(x))(N)F' \rightarrow_{\beta} E\lambda x(M[x/z])(N)F' \rightarrow_{\beta} EM[x/z][N/x]F' = EM[N/z]F',$$

where  $EM[x/z][N/x]F' = EM[N/z]F'$  follows from the fact that  $x$  does not occur free in  $M$ . This is a sequence of  $\beta$ -reductions followed by an (empty) sequence of  $\eta$ -reductions, as desired.

In case (b), because  $M$  contains  $t$  as a subexpression, we can write  $M$  as  $G[t/u]$ , where  $u$  is a completely new variable. Our reduction sequence then has the form

$$E\lambda z(G[t'/u])(N)F' \rightarrow_{\eta} E\lambda z(G[t/u])(N)F' \rightarrow_{\beta} EG[t/u][N/z]F'$$

which can be rewritten

$$E\lambda z(G[t'/u])(N)F' \rightarrow_{\beta} EG[t'/u][N/z]F' \rightarrow_{\beta} EG[t/u][N/z]F'$$

where we have used Lemma 2.23 to get  $G[t'/u][N/z] \rightarrow_{\beta} G[t/u][N/z]$ , using the fact that  $G[t'/u] \rightarrow_{\beta} G[t/u]$ .

In case (c), using similar considerations our reduction sequence can be written in the form

$$E\lambda z(M)(G[t'/u])F' \rightarrow_{\eta} E\lambda z(M)(G[t/u])F' \rightarrow_{\beta} EM[G[t/u]/z]F'$$

which can be rewritten

$$E\lambda z(M)(G[t'/u])F' \rightarrow_{\beta} EM[G[t'/u]/z]F' \rightarrow_{\eta} EM[G[t/u]/z]F'$$

This completes the argument for case (iii) in which  $s$  has the form  $\lambda z(M)(N)$ .

Next, suppose in case (ii) that  $s$  has the form  $\pi_1\langle M, N \rangle$ . Because  $t$  is a proper subexpression of  $s$ , we have that (a)  $t$  is  $\langle M, N \rangle$ , (b)  $t$  is a subexpression of  $M$ , or (c)  $t$  is a subexpression of  $N$ .

In case (a) we again invoke the fact that we are working in the typed setting. If  $t$  is  $\langle M, N \rangle$ , then the type of  $t$  has the form  $\alpha \times \beta$ . The only  $\eta$ -reduction that returns an object of such a type is the reduction rule  $\langle \pi_1(S), \pi_2(S) \rangle \rightarrow_{\eta} S$ . Thus,  $t'$  has the form  $\langle \pi_1\langle M, N \rangle, \pi_2\langle M, N \rangle \rangle$  and the reduction sequence has the form

$$E \pi_1\langle \pi_1\langle M, N \rangle, \pi_2\langle M, N \rangle \rangle F' \rightarrow_{\eta} E\pi_1\langle M, N \rangle F' \rightarrow_{\beta} EMF'$$

But then we also have the reduction sequence

$$E \pi_1\langle \pi_1\langle M, N \rangle, \pi_2\langle M, N \rangle \rangle F' \rightarrow_{\beta} E\pi_1\langle M, N \rangle F' \rightarrow_{\beta} EMF'$$

This is a sequence of  $\beta$ -reductions followed by an (empty) sequence of  $\eta$ -reductions, as desired.

In case (b), our reduction sequence can be written in the form

$$E' \pi_1\langle G[t/u], N \rangle F' \rightarrow_{\eta} E' \pi_1\langle G[t/u], N \rangle F' \rightarrow_{\beta} E'G[t/u]F'$$

which may be rewritten

$$E' \pi_1\langle G[t/u], N \rangle F' \rightarrow_{\beta} E'G[t/u]F' \rightarrow_{\eta} E'G[t/u]F'.$$

Case (c) is completely symmetric. The case in which  $s$  has the form  $\pi_2\langle M, N \rangle$  is also completely symmetric.

Next, suppose in case (ii) that  $s$  has the form  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)$ . Because  $t$  is a proper subexpression of  $s$ , we have that (a)  $t$  is the subexpression  $\mathbf{in}_1(R)$ , (b)  $t$  is a subexpression of  $R$ , (c)  $t$  is a subexpression of  $P$ , or (d)  $t$  is a subexpression of  $Q$ .

In case (a), we again use some facts about typing. If  $t$  is  $\mathbf{in}_1(R)$ , then the type of  $t$  must have the form  $\alpha + \beta$ . The only  $\eta$ -reduction that returns an object of such a type is the reduction rule  $\mathbf{case}(t, [x]\mathbf{in}_1(x), [y]\mathbf{in}_2(y)) \rightarrow_\eta t$ , and thus we know that  $t'$  has the form  $\mathbf{case}(\mathbf{in}_1(R), [x]\mathbf{in}_1(x), [y]\mathbf{in}_2(y))$ , which both  $\eta$  and  $\beta$  reduces to  $\mathbf{in}_1(R)$ . The reduction sequence thus has the form

$$\begin{aligned} E' \mathbf{case}(\mathbf{case}(\mathbf{in}_1(R), [x]\mathbf{in}_1(x), [y]\mathbf{in}_2(y)), [u]P, [v]Q)F' \\ \rightarrow_\eta E' \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)F' \rightarrow_\beta E' P[R/u]F' \end{aligned}$$

which may be rewritten

$$\begin{aligned} E' \mathbf{case}(\mathbf{case}(\mathbf{in}_1(R), [x]\mathbf{in}_1(x), [y]\mathbf{in}_2(y)), [u]P, [v]Q)F' \\ \rightarrow_\beta E' \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)F' \rightarrow_\beta E' P[R/u]F'. \end{aligned}$$

This is a sequence of  $\beta$ -reductions followed by an (empty) sequence of  $\eta$ -reductions, as desired.

In case (b), we may rewrite  $R$  as  $G[t/w]$  (where  $w$  is a new variable), and the reduction sequence has the form

$$E' \mathbf{case}(\mathbf{in}_1(G[t'/w]), [u]P, [v]Q)F' \rightarrow_\eta E' \mathbf{case}(\mathbf{in}_1(G[t/w]), [u]P, [v]Q)F' \rightarrow_\beta E' P[G[t/w]/u]F'$$

which can be rewritten

$$E' \mathbf{case}(\mathbf{in}_1(G[t'/w]), [u]P, [v]Q)F' \rightarrow_\beta E' P[G[t'/w]/u]F' \rightarrow_\eta E' P[G[t/w]/u]F'$$

In case (c), our reduction sequence then has the form

$$E' \mathbf{case}(\mathbf{in}_1(R), [u]G[t'/w], [v]Q)F' \rightarrow_\eta E' \mathbf{case}(\mathbf{in}_1(R), [u]G[t/w], [v]Q)F' \rightarrow_\beta E' G[t/w][R/u]F'$$

which may be rewritten

$$E' \mathbf{case}(\mathbf{in}_1(R), [u]G[t'/w], [v]Q)F' \rightarrow_\beta E' G[t'/w][R/u]F' \rightarrow_\eta E' G[t/w][R/u]F'$$

where again we have used Lemma 2.23 to get  $G[t'/w][R/u] \rightarrow_\eta G[t/w][R/u]$ , using the fact that  $G[t'/w] \rightarrow_\eta G[t/w]$ .

In case (d) our reduction sequence has the form

$$E' \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]G[t'/w])F' \rightarrow_\eta E' \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]G[t/w])F' \rightarrow_\beta E' P[R/u]F'$$



which may be rewritten

$$E' \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]G[t'/w])F' \rightarrow_{\beta} E' P[R/u]F'.$$

This completes the argument for case (iii) in which  $s$  has the form  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)$ . The final case in which the redex  $s$  has the form  $\mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q)$  is completely symmetric.

**Theorem 2.24: The Church-Rosser Theorem for  $\beta$ - and  $\beta\eta$ -reductions in  $\lambda_1$ .**

In the untyped  $\lambda$ -calculus  $\lambda_1$ , if  $X \rightarrow_{\beta} Y_1$  and  $X \rightarrow_{\beta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ . Likewise, if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

**Proof of Theorem 2.24.**

In this proof, all  $\lambda$ -terms are understood to be terms of the untyped  $\lambda$ -calculus  $\lambda_0$ . As before, we assume all terms obey the Barendregt variable convention. The proof follows the same strategy as the proofs of Theorems 1.13 and 1.20, with which familiarity is assumed.

For ease of reference, we list all redexes and their corresponding rules of  $\beta\eta$ -reduction together:

- (i) the redex  $\lambda x(M)(N)$  reduces to  $M[N/x]$
- (ii) the redex  $\pi_1\langle M, N \rangle$  reduces to  $M$
- (iii) the redex  $\pi_2\langle M, N \rangle$  reduces to  $N$
- (iv) the redex  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)$  reduces to  $P[R/u]$
- (v) the redex  $\mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q)$  reduces to  $Q[R/v]$
- (vi) the redex  $\lambda x(M(x))$  reduces to  $M$  where  $x$  is not free in  $M$
- (vii) the redex  $\langle \pi_1(M), \pi_2(M) \rangle$  reduces to  $M$
- (viii) the redex  $\mathbf{case}(M, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v))$  reduces to  $M$

We use the term ‘type (n) redex’ to refer to the redexes in this list. (For example, a type (vii) redex is one of the form  $\langle \pi_1(M), \pi_2(M) \rangle$ .)

As before, we define a more general notion of  $\beta\eta$ -reduction that we denote  $\Rightarrow_{\beta\eta}$ . To begin, note that as before for a given  $\lambda$ -term  $e$  the set of  $\beta$ - or  $\eta$ -redexes contained in  $e$  form a finite set of trees when ordered under the relation of inclusion (with the innermost redexes located at the uppermost leaves of these trees.) We would like it to be the case that whenever one redex  $r$  lies higher than another redex  $r'$  in one of these trees, then when the redex  $r'$  is written in the appropriate form (i)-(viii), we have that  $r$  is a subexpression of  $M, N, P, Q$  or  $R$ . This guarantees that if  $r$  is reduced,  $r'$  still maintains its original form as one of our redex types.

As in the proof of Theorem 1.20, the problem is that there are redexes  $r'$  containing redexes  $r$  as subexpressions such that  $r$  itself is not necessarily a subexpression of the relevant  $M, N, P, Q$  or  $R$ . To help us systematically identify all such cases, we begin by enumerating all subexpressions of redexes that are not subexpressions of the relevant  $M, N, P, Q$  or  $R$ .

- (1)  $\lambda x(M)$  is a proper subexpression of a type (i) redex
- (2)  $\langle M, N \rangle$  is a proper subexpression of a type (ii) or (iii) redex
- (3)  $\mathbf{in}_1(R)$  and  $\mathbf{in}_2(R)$  are proper subexpressions of type (iv) and (v) redexes
- (4)  $M(x)$  is a proper subexpression of a type (vi) redex
- (5)  $\pi_1(M)$  and  $\pi_2(M)$  are proper subexpression of a type (vii) redex
- (6)  $\mathbf{in}_1(u)$  and  $\mathbf{in}_2(u)$  are proper subexpression of a type (viii) redex

We consider each of these cases in turn. In case (1), the only way  $\lambda x(M)$  can be a redex is if it is a type (vi) redex. In case (2), the only way a term  $\langle M, N \rangle$  can be a redex is if it is a type (vii) redex. In case (3), the terms  $\mathbf{in}_1(R)$  and  $\mathbf{in}_2(R)$  never have the form of a redex. In case (4), the only way a term  $M(x)$  can be a redex is if it is a type (i) redex. In case (5), the only way  $\pi_1(M)$  or  $\pi_2(M)$  can be a redex is if it is a type (ii) or (iii) redex. And in case (6),  $\mathbf{in}_1(u)$  and  $\mathbf{in}_2(u)$  can never have the form of a redex.

Thus, the redexes and corresponding subredexes we need to pay special attention to are the following, listed with their types and with the problematic subredexes explicitly indicated

- |   |  |
|---|--|
| <b>(a)</b> the type (i) redex           | $\underbrace{\lambda x(\overline{M}(x))}_{\eta\text{-redex}}(N)$   |
| <b>(b)</b> the type (ii) or (iii) redex | $\pi_i \underbrace{\langle \pi_1(M), \pi_2(M) \rangle}_{\eta\text{-redex}}$                                |
| <b>(c)</b> the type (vi) redex          | $\lambda x \underbrace{(\lambda z(\overline{M})(x))}_{\beta\text{-redex}}$                                 |
| <b>(d)</b> the type (vii) redex         | $\underbrace{\langle \pi_1 \langle M, N \rangle, \pi_2 \langle M, N \rangle \rangle}_{\beta\text{-redex}}$ |

In **(a)**, it is assumed that  $x$  is not free in  $\overline{M}$ , and in **(c)** it is assumed that  $x$  is not free in  $\lambda z(\overline{M})$ . These redexes correspond to redex trees that contain the following structures

$$\begin{array}{cccc}
 \lambda x(\overline{M}(x)) & \langle \pi_1(M), \pi_2(M) \rangle & \lambda x(\overline{M}(x)) & \pi_1 \langle M, N \rangle \quad \pi_2 \langle M, N \rangle \\
 | & | & | & \diagdown \quad \diagup \\
 \lambda x(\overline{M}(x))(N) & \pi_i \langle \pi_1(M), \pi_2(M) \rangle & \lambda x(\overline{M}(x))(N) & \langle \pi_1 \langle M, N \rangle, \pi_2 \langle M, N \rangle \rangle
 \end{array}$$

In analogy with the proof of Theorem 1.20, we declare redexes of the form **(a)-(d)** *prohibited*. If we perform a 1-step reduction on any non-prohibited redex, then any other non-prohibited redex below it maintains its form as a redex. We can therefore perform a 1-step  $\beta$ - or  $\eta$ -reduction on as many (or as few) of the non-prohibited redexes contained in the redex trees

of  $e$  that we want, working from the innermost redexes outwards – that is, working from the uppermost redexes downwards. As before, if doing so produces a term  $e'$ , we will say that  $e \Rightarrow_{\beta} e'$ . If in this way we reduce *all* the non-prohibited redexes contained in  $e$ , we call the resulting term  $e^*$  (or  $[e]^*$ ), which we call the *complete reduction* of  $e$ . (Again, the order in which we perform  $\beta$ -reductions associated with disjoint redexes does not matter.) By definition, we then have  $e \Rightarrow_{\beta} e^*$ .

We now have a sequence of lemmas to those contained in the proofs of Theorem 1.13 and 1.20. Recall that  $s \rightarrow_{\beta\eta} t$  means  $s \rightarrow_{\beta} t$  or  $s \rightarrow_{\eta} t$ .

**Lemma 1:** For any  $\lambda$ -terms  $s$  and  $t$ ,  $s \rightarrow_{\beta\eta} t$  entails  $s \Rightarrow_{\beta\eta} t$ .

**Proof:** The reduction  $s \rightarrow_{\beta\eta} t$  corresponds to the case in which a single redex  $r$  in the redex trees of  $s$  is reduced, with the result being  $t$ . If the redex  $r$  is not prohibited, then it is immediate that  $s \Rightarrow_{\beta\eta} t$ .

In order to deal with the case in which  $r$  is prohibited, we argue for the following slightly modified version of (23) used in the proof of Theorem 1.20

$$\text{If } s \rightarrow_{\beta\eta} t \text{ via the reduction of a prohibited redex } r, \text{ then there are disjoint redexes } r'_1, \dots, r'_n \text{ properly contained in } r \text{ such that } s \rightarrow_{\beta\eta} t \text{ via the reduction of the redexes } r'_1, \dots, r'_n. \quad (56)$$

To argue for this principle, it suffices to consider the four possible forms **(a)**, **(b)**, **(c)** or **(d)** that the prohibited redex  $r$  might have. In the case in which  $r$  has the form **(a)** or **(c)**, the analysis given in the discussion of (23) in the proof of Theorem 1.20 shows that (56) holds.

In case **(b)**,  $r$  has the form of the  $\beta$ -redex  $\pi_i \langle \pi_1(M), \pi_2(M) \rangle$ . This  $\beta$ -reduces to  $\pi_i(M)$ . Note however that if in  $r$  the  $\eta$ -redex  $r' = \langle \pi_1(M), \pi_2(M) \rangle$  is  $\eta$ -reduced, the result is also  $\pi_i(M)$ , and so (56) holds.

In case **(d)**,  $r$  has the form of the  $\eta$ -redex  $\langle \pi_1 \langle M, N \rangle, \pi_2 \langle M, N \rangle \rangle$ . This  $\eta$ -reduces to  $\langle M, N \rangle$ . Note however that if in  $r$  the disjoint  $\beta$ -redexes  $r'_1 = \pi_1 \langle M, N \rangle$  and  $r'_2 = \pi_2 \langle M, N \rangle$  are  $\beta$ -reduced, the result is also  $\langle M, N \rangle$ , and so (56) holds.

Thus, (56) holds in general. Of course, any one of the redexes  $r'_1, \dots, r'_n$  given in (56) may also be prohibited. In this case, the principle (56) can be applied again to any such prohibited  $r'_i$ . Because we cannot have an infinite sequence  $u, u', u'', \dots$  of redexes, each a subredex of the previous one on the list, it follows that if  $s \rightarrow_{\beta\eta} t$  via the reduction of a prohibited redex  $r$ , then there must be a set of disjoint, non-prohibited redexes  $r^*_1, \dots, r^*_m$  properly contained in  $r$  such that  $s \rightarrow_{\beta\eta} t$  via the reduction of the redexes  $r^*_1, \dots, r^*_m$ . Thus,  $s \Rightarrow_{\beta\eta} t$ .  $\square$

**Lemma 2:** For any  $\lambda$ -terms  $s$  and  $t$ ,  $s \Rightarrow_{\beta\eta} t$  entails  $s \rightarrow_{\beta\eta} t$ .

**Proof:** Identical to that given in the proof of Theorem 1.13  $\square$

**Lemma 3:** For any  $\lambda$ -terms  $s, s', t, t'$ , if  $s \Rightarrow_{\beta\eta} s'$  and  $t \Rightarrow_{\beta\eta} t'$ , then  $s[t/v] \Rightarrow_{\beta\eta} s'[t'/v]$ .

**Proof:** Identical to that given in the proof of Theorem 1.20  $\square$

**Lemma 4:** For any  $\lambda$ -terms  $s$  and  $t$ , if  $s \Rightarrow_{\beta\eta} t$  then  $t \Rightarrow_{\beta\eta} s^*$ .

**Proof:** As in the proof of Theorem 1.20, we use the notation  $W(X, Y, \dots)$  to compactly denote

redexes and  $W^{\mathbf{r}}(X, Y, \dots)$  their reductions as follows

type <b>(i)</b> redex:	$W(M, N) = \lambda x(M)(N)$	$W^{\mathbf{r}}(M, N) = M[N/x]$
type <b>(ii)</b> redex:	$W(M, N) = \pi_1\langle M, N \rangle$	$W^{\mathbf{r}}(M, N) = M$
type <b>(iii)</b> redex:	$W(M, N) = \pi_2\langle M, N \rangle$	$W^{\mathbf{r}}(M, N) = N$
type <b>(iv)</b> redex:	$W(P, Q, R) = \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)$	$W^{\mathbf{r}}(P, Q, R) = P[R/u]$
type <b>(v)</b> redex:	$W(P, Q, R) = \mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q)$	$W^{\mathbf{r}}(P, Q, R) = Q[R/v]$
type <b>(vi)</b> redex:	$W(M) = \lambda x(M(x))$	$W^{\mathbf{r}}(M, N) = M$
type <b>(vii)</b> redex:	$W(M) = \langle \pi_1(M), \pi_2(M) \rangle$	$W^{\mathbf{r}}(M, N) = M$
type <b>(viii)</b> redex:	$W(M) = \mathbf{case}(M, [u]\mathbf{in}_1(u), [v]\mathbf{in}_2(v))$	$W^{\mathbf{r}}(M, N) = M$

in type **(vi)** redexes,  $x$  is of course assumed not to be free in  $M$ .

First we show

**(P1)** If a redex  $r = W(X, Y, \dots)$  is not prohibited and another redex  $r'$  is properly contained in  $r$ , then  $r'$  is properly contained in one of  $X, Y, \dots$

**(P2)** For any redex  $W(X, Y, \dots)$ , if  $X \Rightarrow_{\beta\eta} X', Y \Rightarrow_{\beta\eta} Y', \dots$  then

$$W^{\mathbf{r}}(X, Y, \dots) \Rightarrow_{\beta\eta} W^{\mathbf{r}}(X', Y', \dots).$$

The property **(P1)** follows from our characterization of prohibited redexes. For a type **(i)** redex, property **(P2)** amounts to the claim that if  $M \Rightarrow_{\beta\eta} M'$  and  $N \Rightarrow_{\beta\eta} N'$ , then  $M[N/x] \Rightarrow_{\beta\eta} M'[N'/x]$ , which follows from Lemma 3. Likewise for type **(iv)** and **(v)** redexes. For a type **(ii)** redex, property **(P2)** amounts to the trivial claim that if  $M \Rightarrow_{\beta\eta} M'$  and  $N \Rightarrow_{\beta\eta} N'$ , then  $M \Rightarrow_{\beta\eta} M'$ . Likewise for type **(iii)**, **(vi)** **(vii)**, and **(viii)** redexes. So both **(P1)** and **(P2)** are true in general. Having established these, the remainder of the proof of Lemma 4 is identical to that given in Theorem 1.20.  $\square$

Given Lemma 4, the proof of Lemma 5 and remainder of the proof is the same as before, with only trivial notational changes.  $\square$

### Theorem 2.26: The Church-Rosser Theorem for $\beta$ - and $\beta\eta$ -reductions in $\lambda_1^{\text{type}}$ .

In the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$ , if  $X \rightarrow_{\beta} Y_1$  and  $X \rightarrow_{\beta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta} Z$  and  $Y_2 \rightarrow_{\beta} Z$ . Likewise, if  $X \rightarrow_{\beta\eta} Y_1$  and  $X \rightarrow_{\beta\eta} Y_2$ , then there is a term  $Z$  such that  $Y_1 \rightarrow_{\beta\eta} Z$  and  $Y_2 \rightarrow_{\beta\eta} Z$ .

#### Proof of Theorem 2.26.

To prove Theorems 2.26, it suffices to use the proof of Theorems 2.24 and and note that noth-

ing more than notational alterations are required in the typed case. In particular, every term used in these proofs can be typed appropriately if required.  $\square$

**Theorem 2.28: Strong Normalizability of  $\lambda_1^{\text{type}}$ .**

There is no infinite sequence  $\tau_1, \tau_2, \dots$  of terms in the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$  for which

$$\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$$

**Proof of Theorem 2.28.**

This proof is adapted from Wojdyga [11]. We show that if there were an infinite sequence of terms  $\tau_1, \tau_2, \dots$  in the typed  $\lambda$ -calculus  $\lambda_1^{\text{type}}$  for which  $\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$ , then there would also be an infinite sequence of terms  $\tau'_1, \tau'_2, \dots$  in the typed  $\lambda$ -calculus  $\lambda_0^{\text{type}}$  for which  $\tau'_1 \rightarrow_{\beta\eta} \tau'_2 \rightarrow_{\beta\eta} \tau'_3 \rightarrow_{\beta\eta} \dots$ , thereby contradicting the strong normalizability result for  $\lambda_0^{\text{type}}$  proved in Theorem 1.33.

Pick a fundamental simple type of  $\lambda_0^{\text{type}}$  and denote it  $\perp$ . To each type  $P$  of  $\lambda_1^{\text{type}}$  (i.e., each propositional type) we associate a type  $\|P\|$  of  $\lambda_0^{\text{type}}$  (i.e., a simple type) by applying the following rules:

- $\|P\| = \perp$  if  $M$  is a fundamental type or 0
- $\|P_1 \rightarrow P_2\| = \|P_1\| \rightarrow \|P_2\|$
- $\|P_1 \times P_2\| = (\|P_1\| \rightarrow (\|P_2\| \rightarrow \perp)) \rightarrow \perp$
- $\|P_1 + P_2\| = (\|P_1\| \rightarrow \perp) \rightarrow ((\|P_2\| \rightarrow \perp) \rightarrow \perp)$

We then have

**Lemma 1:** For each type  $P$  of  $\lambda_1^{\text{type}}$ ,  $\|P\|$  can be written uniquely in the form

$$P_1 \rightarrow (P_2 \rightarrow (\dots \rightarrow (P_n \rightarrow \perp) \dots)).$$

**Proof:** A straightforward induction on the construction of  $P$ .  $\square$

To each  $\lambda$ -term  $\tau$  of  $\lambda_1^{\text{type}}$ , we associate a  $\lambda$ -term  $\|\tau\|$  of  $\lambda_0^{\text{type}}$  by applying the following rules

- (1)  $\|v^{\tau}\| = v^{\|\tau\|}$
- (2)  $\|M^{\alpha \rightarrow \beta}(N^{\alpha})\| = \|M^{\alpha \rightarrow \beta}\|(\|N^{\alpha}\|)$
- (3)  $\|\lambda x^{\alpha}(M^{\beta})\| = \lambda x^{\|\alpha\|}(\|M^{\beta}\|)$
- (4)  $\|\langle M^{\alpha}, N^{\beta} \rangle\| = \lambda z^{\|\alpha\| \rightarrow (\|\beta\| \rightarrow \perp)}(z(\|M^{\alpha}\|)(\|N^{\beta}\|))$
- (5)  $\|\pi_1(M^{\alpha \times \beta})\| = \lambda z_1^{\alpha_1} \dots \lambda z_n^{\alpha_n}(\|M^{\alpha \times \beta}\|(\lambda u^{\|\alpha\|} \lambda v^{\|\beta\|}(uz_1 \dots z_n)))$   
where  $\|\alpha\| = \alpha_1 \rightarrow (\alpha_2 \rightarrow (\dots \rightarrow (\alpha_n \rightarrow \perp) \dots))$ .

- (6)  $\|\pi_2(M^{\alpha \times \beta})\| = \lambda z_1^{\beta_1} \dots \lambda z_n^{\beta_n} (\|M^{\alpha \times \beta}\| (\lambda u^{\|\alpha\|} \lambda v^{\|\beta\|} (v z_1 \dots z_n)))$   
 where  $\|\beta\| = \beta_1 \rightarrow (\beta_2 \rightarrow (\dots \rightarrow (\beta_n \rightarrow \perp) \dots))$ .
- (7)  $\|\mathbf{in}_1^{\alpha, \beta}(M^\alpha)\| = \lambda u^{\|\alpha\| \rightarrow \perp} \lambda v^{\|\beta\| \rightarrow \perp} (u \|M^\alpha\|)$
- (8)  $\|\mathbf{in}_2^{\alpha, \beta}(N^\beta)\| = \lambda u^{\|\alpha\| \rightarrow \perp} \lambda v^{\|\beta\| \rightarrow \perp} (v \|N^\beta\|)$
- (9)  $\|\mathbf{case}(R^{\alpha + \beta}, [u^\alpha][P]^\gamma, [v^\beta][Q]^\gamma)\| =$   
 $\lambda z_1^{\gamma_1} \dots \lambda z_n^{\gamma_n} (\|R^{\alpha + \beta}\| (\lambda u^{\|\alpha\|} (\|P^\gamma\| z_1 \dots z_n)) (\lambda v^{\|\beta\|} (\|Q^\gamma\| z_1 \dots z_n)))$   
 where  $\|\gamma\| = \gamma_1 \rightarrow (\gamma_2 \rightarrow (\dots \rightarrow (\gamma_n \rightarrow \perp) \dots))$ .
- (10)  $\|\mathbf{abs}^\alpha([M]^0)\| = \lambda z_1^{\alpha_1} \dots \lambda z_n^{\alpha_n} (\|M^0\|)$   
 where  $\|\alpha\| = \alpha_1 \rightarrow (\alpha_2 \rightarrow (\dots \rightarrow (\alpha_n \rightarrow \perp) \dots))$ .

where the bound variables  $z, z_1, \dots, z_n, u, v$  introduced on the right hand side of these expressions are assumed to not appear anywhere on the left hand side. (Again, we do not write all the superscripts in these typed  $\lambda$ -terms, but include enough superscripts that the type of any subexpression may be easily determined.)

We then have

**Lemma 2:** For each  $\lambda$ -term  $P$  of  $\lambda_1^{\text{type}}$  of type  $\tau$ ,  $\|P\|$  is a term of  $\lambda_0^{\text{type}}$  of type  $\|\tau\|$ .

**Proof:** By induction on the construction of  $P$ . There are 10 cases to consider, corresponding to each of the 10 cases in the definition of  $\|\tau\|$  just given. The arguments are all straightforward.

For an easy example, in case (2) we must show that  $\|M^{\alpha \rightarrow \beta}(N^\alpha)\|$  is a  $\lambda_0^{\text{type}}$  term of type  $\|\beta\|$ . By inductive hypothesis  $\|M^{\alpha \rightarrow \beta}\|$  and  $\|N^\alpha\|$  are  $\lambda_0^{\text{type}}$  terms of type  $\|\alpha \rightarrow \beta\|$  and  $\|\beta\|$  respectively. The type  $\|\alpha \rightarrow \beta\|$  is by definition  $\|\alpha\| \rightarrow \|\beta\|$ . It follows immediately that  $\|M^{\alpha \rightarrow \beta}(N^\alpha)\| = \|M^{\alpha \rightarrow \beta}\| (\|N^\alpha\|)$  is a  $\lambda_0^{\text{type}}$  term of type  $\|\beta\|$ .

For a harder example, consider case (5). We must show that  $\|\pi_1(M^{\alpha \times \beta})\|$  is a  $\lambda_0^{\text{type}}$  term of type  $\|\alpha\|$ . If  $\|\alpha\| = \alpha_1 \rightarrow (\alpha_2 \rightarrow (\dots \rightarrow (\alpha_n \rightarrow \perp) \dots))$ , then if  $u$  has type  $\|\alpha\|$  and  $z_1, \dots, z_n$  have types  $\alpha_1, \dots, \alpha_n$ , then  $u z_1 \dots z_n$  is a term of type  $\perp$ . Thus  $\lambda u^{\|\alpha\|} \lambda v^{\|\beta\|} (u z_1 \dots z_n)$  is a term of type  $\|\alpha\| \rightarrow (\|\beta\| \rightarrow \perp)$ . By inductive hypothesis,  $\|M^{\alpha \times \beta}\|$  has type  $\|\alpha \times \beta\|$  which is by definition  $(\|\alpha\| \rightarrow (\|\beta\| \rightarrow \perp)) \rightarrow \perp$ , and so  $\|M^{\alpha \times \beta}\| (\lambda u^{\|\alpha\|} \lambda v^{\|\beta\|} (u z_1 \dots z_n))$  has type  $\perp$ . Because  $z_1, \dots, z_n$  have types  $\alpha_1, \dots, \alpha_n$ , it then follows that  $\lambda z_1^{\alpha_1} \dots \lambda z_n^{\alpha_n} (\|M^{\alpha \times \beta}\| (\lambda u^{\|\alpha\|} \lambda v^{\|\beta\|} (u z_1 \dots z_n)))$  has type  $\alpha_1 \rightarrow (\alpha_2 \rightarrow (\dots \rightarrow (\alpha_n \rightarrow \perp) \dots))$ , i.e., type  $\|\alpha\|$ . Thus  $\|\pi_1(M^{\alpha \times \beta})\|$  is a  $\lambda_0^{\text{type}}$  term of type  $\|\alpha\|$ .

Straightforward reasoning of this sort suffices for verifying the remaining 8 cases.  $\square$

The following lemma also holds

**Lemma 3:** If  $M^\tau, N^\sigma$  are  $\lambda_1^{\text{type}}$  terms and  $x^\sigma$  a variable, then  $\|M^\tau\| [\|N^\sigma\| / x^{\|\sigma\|}] = \|M^\tau[N^\sigma / x^\sigma]\|$ .

**Proof:** Straightforward induction on the construction of the  $\lambda$ -term  $M^\tau$ .  $\square$

If  $\sigma$  and  $\tau$  are  $\lambda$ -terms of  $\lambda_0^{\text{type}}$ , we write  $\sigma \rightarrow_{\beta\eta}^+ \tau$  just in case there is a  $\beta\eta$ -reduction of  $\sigma$  to  $\tau$  involving at least one step (i.e., there is one such reduction which is not a trivial 0-step reduction.) We then have

**Lemma 4:** For terms  $M, N$  of  $\lambda_1^{\text{type}}$ , if  $M \rightarrow_\beta N$ , then  $\|M\| \rightarrow_{\beta\eta}^+ \|N\|$ .

**Proof:** We consider each of the 5  $\beta$ -reduction rules of  $\lambda_1^{\text{type}}$  in turn. First, consider the reduction rule

$$\lambda x^\sigma (M^\tau)(N^\sigma) \rightarrow_\beta M^\tau [N^\sigma / x^\sigma].$$

Then

$$\begin{aligned} \|\lambda x^\sigma (M^\tau)(N^\sigma)\| &= \|\lambda x^\sigma (M^\tau)\|(\|N^\sigma\|) = \lambda x^{\|\sigma\|}(\|M^\tau\|)(\|N^\sigma\|) \\ &\rightarrow_\beta \|M^\tau\|[\|N^\sigma\|/x] = \|M^\tau[N^\sigma/x]\| \end{aligned}$$

as desired. Next consider the reduction rule

$$\pi_1(\langle M^\sigma, N^\tau \rangle^{\sigma \times \tau}) \rightarrow_\beta M^\sigma$$

Then

$$\begin{aligned} \|\pi_1(\langle M^\sigma, N^\tau \rangle^{\sigma \times \tau})\| &= \lambda z_1^{\sigma_1} \dots \lambda z_n^{\sigma_n} (\|\langle M^\sigma, N^\tau \rangle\|(\lambda u^{\|\sigma\|} \lambda v^{\|\tau\|}(uz_1 \dots z_n))) \\ &= \lambda z_1^{\sigma_1} \dots \lambda z_n^{\sigma_n} (\lambda z^{\|\sigma\| \rightarrow (\|\tau\| \rightarrow \perp)}(z(\|M^\sigma\|)(\|N^\tau\|))(\lambda u^{\|\sigma\|} \lambda v^{\|\tau\|}(uz_1 \dots z_n))) \\ &\rightarrow_\beta \lambda z_1^{\sigma_1} \dots \lambda z_n^{\sigma_n} (\lambda u^{\|\sigma\|} \lambda v^{\|\tau\|}(uz_1 \dots z_n)(\|M^\sigma\|)(\|N^\tau\|)) \\ &\rightarrow_\beta \lambda z_1^{\sigma_1} \dots \lambda z_n^{\sigma_n} (\|M^\sigma\|z_1 \dots z_n) \rightarrow_\eta \lambda z_1^{\sigma_1} \dots \lambda z_{n-1}^{\sigma_{n-1}} (\|M^\sigma\|z_1 \dots z_{n-1}) \rightarrow_\eta \dots \rightarrow_\eta \|M^\sigma\| \end{aligned}$$

where  $\|\sigma\| = \sigma_1 \rightarrow (\sigma_2 \rightarrow (\dots \rightarrow (\sigma_n \rightarrow \perp) \dots))$ , as desired. The argument for the reduction rule

$$\pi_2(\langle M^\sigma, N^\tau \rangle^{\sigma \times \tau}) \rightarrow_\beta N^\tau$$

is symmetric. Next consider the reduction rule

$$\mathbf{case}(\mathbf{in}_1^{\alpha, \beta}(S^\alpha)^{\alpha + \beta}, [u^\alpha]P^\gamma, [v^\beta]Q^\gamma) \rightarrow_\beta P^\gamma[S^\alpha/u^\alpha]$$

Then

$$\begin{aligned} \|\mathbf{case}(\mathbf{in}_1^{\alpha, \beta}(S^\alpha)^{\alpha + \beta}, [u^\alpha]P^\gamma, [v^\beta]Q^\gamma)\| &= \lambda z_1^{\gamma_1} \dots \lambda z_n^{\gamma_n} (\|\mathbf{in}_1^{\alpha, \beta}(S^\alpha)^{\alpha + \beta}\|(\lambda u^{\|\alpha\|}(\|P^\gamma\|z_1 \dots z_n))(\lambda v^{\|\beta\|}(\|Q^\gamma\|z_1 \dots z_n))) \\ &= \lambda z_1^{\gamma_1} \dots \lambda z_n^{\gamma_n} (\lambda \bar{u}^{\|\alpha\| \rightarrow \perp} \lambda \bar{v}^{\|\beta\| \rightarrow \perp} (\bar{u}\|S^\alpha\|)(\lambda u^{\|\alpha\|}(\|P^\gamma\|z_1 \dots z_n))(\lambda v^{\|\beta\|}(\|Q^\gamma\|z_1 \dots z_n))) \\ &\rightarrow_\beta \lambda z_1^{\gamma_1} \dots \lambda z_n^{\gamma_n} (\lambda u^{\|\alpha\|}(\|P^\gamma\|z_1 \dots z_n)\|S^\alpha\|) \rightarrow_\beta \lambda z_1^{\gamma_1} \dots \lambda z_n^{\gamma_n} (\|P^\gamma\|[\|S^\alpha\|/u^{\|\alpha\|}]z_1 \dots z_n) \\ &\rightarrow_\eta \|P^\gamma\|[\|S^\alpha\|/u^{\|\alpha\|}] = \|P^\gamma[S^\alpha/u^\alpha]\| \end{aligned}$$

where  $\|\gamma\| = \gamma_1 \rightarrow (\gamma_2 \rightarrow (\dots \rightarrow (\gamma_n \rightarrow \perp) \dots))$ , as desired. The argument for the reduction rule

$$\mathbf{case}(\mathbf{in}_2^{\alpha, \beta}(S^\beta)^{\alpha + \beta}, [u^\alpha]P^\gamma, [v^\beta]Q^\gamma) \rightarrow_\beta Q^\gamma[S^\beta/v^\beta]$$

is symmetric.  $\square$

To prove the strong normalizability of  $\lambda_1^{\text{type}}$  under  $\beta$ -reductions, suppose to the contrary that there is an infinite sequence  $\tau_1, \tau_2, \dots$  of simply typed  $\lambda$ -terms such that

$$\tau_1 \rightarrow_{\beta} \tau_2 \rightarrow_{\beta} \tau_3 \rightarrow_{\beta} \dots$$

Then by Lemma 4,

$$\|\tau_1\| \rightarrow_{\beta\eta}^+ \|\tau_2\| \rightarrow_{\beta\eta}^+ \|\tau_3\| \rightarrow_{\beta\eta}^+ \dots$$

contradicting the strong normalizability result for  $\lambda_0^{\text{type}}$  proved in Theorem 1.33.

### Theorem 2.31: Soundness and Completeness of $\mathbf{TR}_1$ .

Suppose  $E$  is an untyped  $\lambda_1$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) It is possible to assign types of  $\lambda_1^{\text{type}}$  to all the variables of  $E$  in such a way that the free variables  $v_1, \dots, v_n$  of  $E$  are assigned types  $\tau_1, \dots, \tau_n$ , and  $E$  itself has type  $X$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_1$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

### Proof of Theorem 2.31.

First we prove (i)  $\rightarrow$  (ii). The proof is by induction on the construction of the  $\lambda$ -term  $E$ . If  $E$  is a variable, or has one of the forms  $E_1(E_2)$  or  $\lambda v(E)$ , then the argument given in the proof of Theorem 1.43 suffices.

Suppose  $E$  has the form  $\langle E_1, E_2 \rangle$ , and that by assigning the free variables  $v_1, \dots, v_n$  of  $\langle E_1, E_2 \rangle$  the types  $\tau_1, \dots, \tau_n$  respectively,  $E$  can be assigned the type  $X \times Y$ . Then by inductive hypothesis we have

$$\Gamma_1 \vdash E_1 : X \quad \text{and} \quad \Gamma_2 \vdash E_2 : Y$$

for some contexts  $\Gamma_1, \Gamma_2$  such that for any  $j$ , if  $v_j$  is free in  $E_1$  (resp.  $E_2$ ), then  $v_j : \tau_j$  appears in  $\Gamma_1$  (resp.  $\Gamma_2$ ). Letting  $\bar{E}_1$  and  $\bar{E}_2$  be  $\lambda$ -terms obeying the Barendregt variable convention and  $\alpha$ -equivalent to  $E_1$  and  $E_2$ , by Theorem 2.30 there are proofs

$$\Gamma'_1 \vdash \bar{E}_1 : X \quad \text{and} \quad \Gamma'_2 \vdash \bar{E}_2 : Y$$

such that  $\Gamma'_1$  (resp.  $\Gamma'_2$ ) is just the set of typing declarations  $v_j : \tau_j$  for the  $v_j$  that are free in  $\bar{E}_1$  (resp.  $\bar{E}_2$ .) Let  $\Gamma$  be the context

$$v_1 : \tau_1, \dots, v_n : \tau_n.$$



By adding additional variable declarations to the contexts of these proofs, we have

$$\Gamma \vdash \bar{E}_1 : X \quad \text{and} \quad \Gamma \vdash \bar{E}_2 : Y.$$

from which we may infer  $\Gamma \vdash \langle E_1, E_2 \rangle : X \times Y$  (where this may involve silent renaming of variables). So **(i)**  $\rightarrow$  **(ii)** holds in this case.

Suppose  $E$  has the form  $\pi_1(E')$ , and that by assigning the free variables  $v_1, \dots, v_n$  of  $\pi_1(E')$  the types  $\tau_1, \dots, \tau_n$  respectively,  $\pi_1(E')$  can be assigned type  $X$ , with  $E'$  being assigned type  $X \times Y$ . Then by inductive hypothesis, we have

$$\Gamma \vdash E' : X \times Y$$

where  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, \dots, v_n : \tau_n$ . We then immediately have

$$\Gamma \vdash \pi_1(E') : X$$

as desired, and so **(i)**  $\rightarrow$  **(ii)** holds in this case.

Similarly straightforward arguments of this sort can be used to deal with the cases in which  $E$  has the form  $\pi_2(E')$ ,  $\mathbf{in}_1(E')$ ,  $\mathbf{in}_2(E')$  or  $\mathbf{abs}(E')$ .

Suppose next that  $E$  has the form  $\mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q)$ . Without loss of generality, we assume that  $E$  obeys the Barendregt variable convention. Assume that by assigning the free variables  $v_1, \dots, v_n$  of  $E$  the types  $\tau_1, \dots, \tau_n$  respectively,  $E$  can be assigned the type  $Z$ . Suppose that in this case  $R$  is assigned type  $X + Y$ ,  $P$  is assigned type  $Z$ , and  $Q$  is assigned type  $Z$ . Then by inductive hypothesis and Theorem 2.30 we have

$$\Gamma_1 \vdash R : X + Y \quad \text{and} \quad \Gamma_2, u : X \vdash P : Z \quad \text{and} \quad \Gamma_3, v : Y \vdash Q : Z$$

where each of  $\Gamma_1, \Gamma_2$  and  $\Gamma_3$  consists of the variable declarations  $v_j : \tau_j$  for any  $v_j$  that appears free in  $R, P$  and  $Q$  respectively. Because  $E$  (and thus all its subterms) obeys the Barendregt variable convention, neither  $u$  or  $v$  appear among the  $v_j$ , and so the contexts  $\Gamma_2, u : X$  and  $\Gamma_3, v : Y$  are grammatical. We thus obtain

$$\Gamma_1, \Gamma_2, \Gamma_3 \vdash \mathbf{case}(\mathbf{in}_1(R), [u]P, [v]Q) : Z$$

as desired. The case in which  $E$  has the form  $\mathbf{case}(\mathbf{in}_2(R), [u]P, [v]Q)$  is of course symmetrical. This completes the argument that **(i)**  $\rightarrow$  **(ii)**.

Now we prove **(ii)**  $\rightarrow$  **(i)**. The proof is by induction on the construction of proofs. If  $\Gamma \vdash E : X$  has the form  $\Gamma', v : X \vdash v : X$ , or is the result of an (App) or (Abs) inference, then the argument given in the proof of Theorem 1.43 suffices. The proofs for the remaining cases are all similar, and we simply consider the most complex case in which the proof is the result of a case inference

$$\frac{\Gamma \vdash R : \sigma + \tau \quad \Gamma, u : \sigma \vdash P : \rho \quad \Gamma, v : \tau \vdash Q : \rho}{\Gamma \vdash \mathbf{case}(R, [u^\sigma]P, [v^\tau]Q) : \rho}$$

Assume then that we have

$$\Gamma \vdash \mathbf{case}(R, [u^\sigma]P, [v^\tau]Q) : \rho,$$

where if  $v_1, \dots, v_n$  are the free variables of  $\mathbf{case}(R, [u^\sigma]P, [v^\tau]Q)$ , then  $\Gamma$  includes typing declarations  $v_1 : \tau_1, \dots, v_n : \tau_n$ . Suppose that the last step in the proof of this sequent is indeed the **case** rule, so that we have

$$\Gamma \vdash R : \sigma + \tau \quad \text{and} \quad \Gamma, u : \sigma \vdash P : \rho \quad \text{and} \quad \Gamma, v : \tau \vdash Q : \rho$$

Then by inductive hypothesis, consistent with the free variables of  $R$  being assigned types in accordance with  $\Gamma$ , the term  $R$  may be assigned the type  $\sigma + \tau$ , consistent with the free variables of  $P$  being assigned types in accordance with  $\Gamma, u : \sigma$ , the term  $P$  may be assigned the type  $\rho$ , and consistent with the free variables of  $Q$  being assigned types in accordance with  $\Gamma, v : \tau$ , the term  $Q$  may be assigned the type  $\rho$ . It follows that consistent with the free variables of  $\mathbf{case}(R, [u^\sigma]P, [v^\tau]Q)$  being assigned types in accordance with  $\Gamma$ , this term may be assigned the type  $\rho$ , as desired. Thus **(ii)**  $\rightarrow$  **(i)** holds.

### Theorem 2.32 Soundness and Completeness of $\mathbf{TR}_1^{\text{bt}}$ .

Suppose  $E$  is a bound-typed  $\lambda_1$ -term whose free variables are  $v_1, \dots, v_n$ , and let  $\tau_1, \dots, \tau_n$  be types. Then the following are equivalent:

- (i) If the variables  $v_1, v_2, \dots, v_n$  are assigned types  $\tau_1, \tau_2, \dots, \tau_n$  respectively, then  $E$  must be assigned type  $\sigma$ .
- (ii) There is a proof of  $\Gamma \vdash E : X$  in  $\mathbf{TR}_1^{\text{bt}}$ , where the context  $\Gamma$  includes the variable declarations  $v_1 : \tau_1, v_2 : \tau_2, \dots, v_n : \tau_n$ .

### Proof of Theorem 2.32.

This proof is largely analogous to that of Theorem 2.31 and so is omitted.

# Bibliography

- [1] H. Barendregt. *The Lambda Calculus, its Syntax and Semantics*. North-Holland, 1984.
- [2] H. Barendregt. The impact of the lambda calculus in logic and computer science. *The Bulletin of Symbolic Logic*, 3(2):181–215, Jun 1997.
- [3] G. Boolos, J. Burgess, and R Jeffrey. *Computability and Logic*. Cambridge University Press, 2007.
- [4] J. Hindley and J. Seldin. *Lambda-Calculus and Combinators, An Introduction*. Cambridge University Press, 2008.
- [5] Hindley J. *Basic Simple Type Theory*. Cambridge University Press, 1997.
- [6] G. Michaelson. *Functional Programming Through Lambda Calculus*. Dover Publications, 1989.
- [7] G. Mints. *A Short Introduction to Intuitionistic Logic*. Kluwer Academic Publishers, 2000.
- [8] B. Pierce. *Types and Programming Languages*. MIT Press, 2002.
- [9] P. Selinger. *Lecture Notes on the Lambda Calculus*. Lulu.com, 2018.
- [10] M. Sorensen and P. Urzyczyn. *Lectures on the Curry-Howard Isomorphism*. Elsevier Science, 2006.
- [11] A. Wojdyga. Short proofs of strong normalization. In E. Ochmański and J. Tyszkiewicz, editors, *Mathematical Foundations of Computer Science 2008. Lecture Notes in Computer Science*, volume 5162, pages 612–623, 2008.