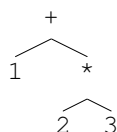


Chapter 1

λ -Calculus

This chapter presents the λ -calculus, a core calculus for functional languages (including SML of course). It captures the essential mechanism of computation in functional languages, and thus serves as an excellent framework for investigating basic concepts in functional languages. According to the Church-Turing thesis, the λ -calculus is equally expressive as Turing machines, but its syntax is deceptively simple. We first discuss the syntax and semantics of the λ -calculus and then show how to write programs in the λ -calculus.

Before we proceed, we briefly discuss the difference between *concrete syntax* and *abstract syntax*. Concrete syntax specifies which string of characters is accepted as a valid program (causing no syntax errors) or rejected as an invalid program (causing syntax errors). For example, according to the concrete syntax of SML, a string `~1` is interpreted as an integer `-1`, but a string `-1` is interpreted as an infix operator `-` applied to an integer argument `1` (which later causes a type error). A parser implementing concrete syntax usually translates source programs into tree structures. For example, a source program `1 + 2 * 3` is translated into



after taking into account operator precedence rules. Such tree structures are called *abstract syntax trees* which abstract away from details of parsing (such as operator precedence/associativity rules) and focus on the structure of source programs; abstract syntax is just the syntax for such tree structures.

While concrete syntax is an integral part of designing a programming language, we will not discuss it in this course. Instead we will work with abstract syntax to concentrate on computational aspects of programming languages. For example, we do not discuss why `1 + 2 * 3` and `1 + (2 * 3)`, both written in concrete syntax, are translated by the parser into the same abstract syntax tree shown above. For the purpose of understanding how their computation proceeds, the abstract syntax tree alone suffices.

1.1 Abstract syntax for the λ -calculus

The abstract syntax for the λ -calculus is given as follows:

$$\text{expression} \quad e ::= x \mid \lambda x. e \mid e e$$

- An expression x is called a *variable*. We may use other names for variables (e.g., $z, s, t, f, arg, accum$, and so on). Strictly speaking, therefore, x itself in the inductive definition of expression is a metavariable.
- An expression $\lambda x. e$ is called a λ -*abstraction*, or just a function, which denotes a mathematical function whose formal argument is x and whose body is e . We may think of $\lambda x. e$ as an internal representation of a nameless SML function $\text{fn } x \Rightarrow e$ in abstract syntax.

We say that a variable x is *bound* in the λ -abstraction $\lambda x. e$ (just like a variable x is bound in an SML function `fn x => e`). Alternatively we can say that x is a bound variable in the body e .

- An expression $e_1 e_2$ is called a λ -*application* or an *application* which denotes a function application (if e_1 is shown to be equivalent to a λ -abstraction somehow). We may think of $e_1 e_2$ as an internal representation of an SML function application in abstract syntax. As in SML, applications are left-associative: $e_1 e_2 e_3$ means $(e_1 e_2) e_3$ instead of $e_1 (e_2 e_3)$.

The scope of a λ -abstraction extends as far to the right as possible. Here are a few examples:

- $\lambda x. x y$ is the same expression as a λ -abstraction $\lambda x. (x y)$ whose body is $x y$. It should not be understood as an application $(\lambda x. x) y$.
- $\lambda x. \lambda y. x y$ is the same expression as $\lambda x. \lambda y. (x y) = \lambda x. (\lambda y. (x y))$. It should not be understood as an application $(\lambda x. \lambda y. x) y$.

As it turns out, every expression in the λ -calculus denotes a mathematical function. That is, the denotation of every expression in the λ -calculus is a mathematical function. Section 1.2 discusses how to determine unique mathematical functions corresponding to expressions in the λ -calculus, and in the present section, we develop the intuition behind the λ -calculus by considering a few examples of λ -abstractions.

Our first example is an identity function:

$$\text{id} = \lambda x. x$$

id is an identity function because when given an argument x , it returns x without any further computation. Like higher-order functions in SML, a λ -abstraction may return another λ -abstraction as its result. For example, tt below takes t to return another λ -abstraction $\lambda f. t$ which ignores its argument; ff below ignores its argument t to return a λ -abstraction $\lambda f. f$:

$$\begin{aligned} \text{tt} &= \lambda t. \lambda f. t = \lambda t. (\lambda f. t) \\ \text{ff} &= \lambda t. \lambda f. f = \lambda t. (\lambda f. f) \end{aligned}$$

Similarly a λ -abstraction may expect another λ -abstraction as its argument. For example, the λ -abstraction below expects another λ -abstraction s which is later applied to z :

$$\text{one} = \lambda s. \lambda z. s z = \lambda s. (\lambda z. (s z))$$

1.2 Operational semantics of the λ -calculus

The semantics of a programming language answers the question of “what is the *meaning* of a given program?” This is an important question in the design of programming languages because lack of formal semantics implies potential ambiguities in interpreting programs. Put another way, lack of formal semantics makes it impossible to determine the meaning of certain programs. Surprisingly not every programming language has its semantics. For example (and perhaps to your surprise), the C language has no formal semantics — the same C program may exhibit different behavior depending on the state of the machine on which the program is executed.

There are three approaches to formulating the semantics of programming languages: *denotational semantics*, *axiomatic semantics*, and *operational semantics*. Throughout this course, we will use exclusively the operational semantics approach for its close connection with judgments and inference rules. The operational semantics approach is also attractive because it directly reflects the implementation of programming languages (e.g., interpreters or compilers).

In general, the operational semantics of a programming language specifies how to transform a program into a *value* via a sequence of “operations.” In the case of the λ -calculus, values consist of λ -abstractions and “operations” are called *reductions*. Thus the operational semantics of the λ -calculus specifies how to reduce an expression e into a value v where v is defined as follows:

$$\text{value} \quad v ::= \lambda x. e$$

Then we take v as the meaning of e . Since a λ -abstraction denotes a mathematical function, it follows that *every* expression in the λ -calculus denotes a mathematical function.

With this idea in mind, let us formally define reductions of expressions. We introduce a *reduction judgment* of the form $e \mapsto e'$.¹

$$e \mapsto e' \quad \Leftrightarrow \quad e \text{ reduces to } e'$$

We write \mapsto^* for the reflexive and transitive closure of \mapsto . That is, $e \mapsto^* e'$ holds if $e \mapsto e_1 \mapsto \dots \mapsto e_n = e'$ where $n \geq 0$. (See Exercise 1.11 for a formal definition of \mapsto^* .) We say that e *evaluates* to v if $e \mapsto^* v$ holds.

Before we provide inference rules to complete the definition of the judgment $e \mapsto e'$, let us see what kind of expression can be reduced to another expression. Clearly variables and λ -abstractions cannot be further reduced:

$$\begin{aligned} x &\not\mapsto \cdot \\ \lambda x. e &\not\mapsto \cdot \end{aligned}$$

($e \not\mapsto \cdot$ means that e does not reduce to another expression.) Then when can we reduce an application $e_1 e_2$? If we think of it as an internal representation of an SML function application, we can reduce it only if e_1 represents an SML function. Thus the only candidate for reduction is an application of the form $(\lambda x. e'_1) e_2$.

If we think of $\lambda x. e'_1$ as a mathematical function whose formal argument is x and whose body is e'_1 , the most natural way to reduce $(\lambda x. e'_1) e_2$ is by substituting e_2 for every occurrence of x in e'_1 , or equivalently, by replacing every occurrence of x in e'_1 by e_2 . (For now, we do not consider the issue of whether e_2 is a value or not.) To this end, we introduce a *substitution* $[e'/x]e$:

$[e'/x]e$ is defined as an expression obtained by substituting e' for every occurrence of x in e .

$[e'/x]e$ may also be read as “applying a substitution $[e'/x]$ to e .” Then the following reduction is justified

$$(\lambda x. e) e' \mapsto [e'/x]e$$

where the expression being reduced, namely $(\lambda x. e) e'$, is called a *redex* (reducible expression). For historical reasons, the above reduction is called a β -reduction.

Simple as it may seem, the precise definition of $[e'/x]e$ is remarkably subtle (see Section 1.3). For now, we just avoid complex examples whose reduction would require the precise definition of substitution. Here are a few examples of β -reductions; the redex for each step is underlined:

$$\begin{aligned} &(\lambda x. x) (\lambda y. y) \mapsto \lambda y. y \\ &\frac{(\lambda t. \lambda f. t) (\lambda x. x) (\lambda y. y)}{(\lambda t. \lambda f. f) (\lambda x. x) (\lambda y. y) \mapsto (\lambda f. \lambda x. x) (\lambda y. y)} \mapsto \lambda x. x \\ &\frac{(\lambda t. \lambda f. f) (\lambda x. x) (\lambda y. y)}{(\lambda s. \lambda z. s z) (\lambda x. x) (\lambda y. y) \mapsto (\lambda f. f) (\lambda y. y)} \mapsto \lambda y. y \\ &\frac{(\lambda s. \lambda z. s z) (\lambda x. x) (\lambda y. y)}{(\lambda s. \lambda z. s z) (\lambda x. x) (\lambda y. y) \mapsto (\lambda z. (\lambda x. x) z) (\lambda y. y)} \mapsto (\lambda x. x) (\lambda y. y) \mapsto \lambda y. y \end{aligned}$$

The β -reduction is the basic principle for reducing expressions, but it does not yield unique inference rules for the judgment $e \mapsto e'$. That is, there can be more than one way to apply the β -reduction to an expression, or equivalently, an expression may contain multiple redexes in it. For example, $(\lambda x. x) ((\lambda y. y) (\lambda z. z))$ contains two redexes in it:

$$\begin{aligned} &\frac{(\lambda x. x) ((\lambda y. y) (\lambda z. z))}{(\lambda x. x) ((\lambda y. y) (\lambda z. z)) \mapsto (\lambda y. y) (\lambda z. z)} \mapsto \lambda z. z \\ &\frac{(\lambda x. x) ((\lambda y. y) (\lambda z. z))}{(\lambda x. x) ((\lambda y. y) (\lambda z. z)) \mapsto (\lambda x. x) (\lambda z. z)} \mapsto \lambda z. z \end{aligned}$$

In the first case, the expression being reduced has the form $(\lambda x. e) e'$ and we immediately apply the β -reduction to the whole expression to obtain $[e'/x]e$. In the second case, we apply the β -reduction to e' which happens to be a redex; if e' was not a redex (e.g., $e' = \lambda t. t$), the second case would be impossible. Here is another example of an expression containing two redexes:

$$\begin{aligned} &\frac{(\lambda s. \lambda z. s z) ((\lambda x. x) (\lambda y. y))}{(\lambda s. \lambda z. s z) ((\lambda x. x) (\lambda y. y)) \mapsto \lambda z. ((\lambda x. x) (\lambda y. y)) z} \mapsto \lambda z. (\lambda y. y) z \mapsto \lambda z. z \\ &\frac{(\lambda s. \lambda z. s z) ((\lambda x. x) (\lambda y. y))}{(\lambda s. \lambda z. s z) ((\lambda x. x) (\lambda y. y)) \mapsto (\lambda s. \lambda z. s z) (\lambda y. y)} \mapsto \lambda z. (\lambda y. y) z \mapsto \lambda z. z \end{aligned}$$

¹After all, the notion of judgment that we learned in Chapter ?? is not really useless!

In the course of reducing an expression to a value, therefore, we may be able to apply the β -reduction in many different ways. As we do not want to apply the β -reduction in an arbitrary way, we need a certain *reduction strategy* so as to apply the β -reduction in a systematic way.

In this course, we consider two reduction strategies: *call-by-name* and *call-by-value*. The call-by-name strategy always reduces the leftmost and outermost redex. To be specific, given an expression $e_1 e_2$, it checks if e_1 is a λ -abstraction $\lambda x. e'_1$. If so, it applies the β -reduction to the whole expression to obtain $[e_2/x]e'_1$. Otherwise it attempts to reduce e_1 using the same reduction strategy without considering e_2 ; when e_1 later reduces to a value (which must be a λ -abstraction), it applies the β -reduction to the whole expression. Consequently the second subexpression in an application (e.g., e_2 in $e_1 e_2$) is never reduced. The call-by-value strategy is similar to the call-by-name strategy, but it reduces the second subexpression in an application to a value v after reducing the first subexpression. Hence the call-by-value strategy applies the β -reduction to an application of the form $(\lambda x. e) v$ only. Note that neither strategy reduces expressions inside a λ -abstraction, which implies that values are not further reduced.

As an example, let us consider an expression $(\text{id}_1 \text{id}_2) (\text{id}_3 (\lambda z. \text{id}_4 z))$ which reduces in different ways under the two reduction strategies; id_i is an abbreviation of an identity function $\lambda x_i. x_i$:

call-by-name	call-by-value
$(\text{id}_1 \text{id}_2) (\text{id}_3 (\lambda z. \text{id}_4 z))$	$(\text{id}_1 \text{id}_2) (\text{id}_3 (\lambda z. \text{id}_4 z))$
$\mapsto \text{id}_2 (\text{id}_3 (\lambda z. \text{id}_4 z))$	$\mapsto \text{id}_2 (\text{id}_3 (\lambda z. \text{id}_4 z))$
$\mapsto \text{id}_3 (\lambda z. \text{id}_4 z)$	$\mapsto \text{id}_2 (\lambda z. \text{id}_4 z)$
$\mapsto \lambda z. \text{id}_4 z$	$\mapsto (\lambda z. \text{id}_4 z)$

The reduction diverges in the second step: the call-by-name strategy applies the β -reduction to the whole expression because it does not need to inspect the second subexpression $\text{id}_3 (\lambda z. \text{id}_4 z)$ whereas the call-by-value strategy chooses to reduce the second subexpression which is not a value yet.

Now we are ready to provide inference rules for the judgment $e \mapsto e'$, which we refer to as *reduction rules*. The call-by-name strategy uses two reduction rules (*Lam* for *Lambda* and *App* for *Application*):

$$\frac{e_1 \mapsto e'_1}{e_1 e_2 \mapsto e'_1 e_2} \text{Lam} \quad \frac{}{(\lambda x. e) e' \mapsto [e'/x]e} \text{App}$$

The call-by-value strategy uses an additional rule to reduce second subexpression in applications; we reuse the reduction rule names from the call-by-name strategy (*Arg* for *Argument*):

$$\frac{e_1 \mapsto e'_1}{e_1 e_2 \mapsto e'_1 e_2} \text{Lam} \quad \frac{e_2 \mapsto e'_2}{(\lambda x. e) e_2 \mapsto (\lambda x. e) e'_2} \text{Arg} \quad \frac{}{(\lambda x. e) v \mapsto [v/x]e} \text{App}$$

A drawback of the call-by-name strategy is that the same expression may be evaluated multiple times. For example, $(\lambda x. x x) ((\lambda y. y) (\lambda z. z))$ evaluates $(\lambda y. y) (\lambda z. z)$ to $\lambda z. z$ eventually twice:

$$\begin{aligned} & (\lambda x. x x) ((\lambda y. y) (\lambda z. z)) \\ \mapsto & ((\lambda y. y) (\lambda z. z)) ((\lambda y. y) (\lambda z. z)) \\ \mapsto & (\lambda z. z) ((\lambda y. y) (\lambda z. z)) \\ \mapsto & (\lambda y. y) (\lambda z. z) \\ \mapsto & \lambda z. z \end{aligned}$$

In the case of the call-by-value strategy, $(\lambda y. y) (\lambda z. z)$ is evaluated only once:

$$\begin{aligned} & (\lambda x. x x) ((\lambda y. y) (\lambda z. z)) \\ \mapsto & (\lambda x. x x) (\lambda z. z) \\ \mapsto & (\lambda z. z) (\lambda z. z) \\ \mapsto & \lambda z. z \end{aligned}$$

On the other hand, the call-by-name strategy never evaluates expressions that do not contribute to evaluations. For example,

$$(\lambda t. \lambda f. f) ((\lambda y. y) (\lambda z. z)) ((\lambda y'. y') (\lambda z'. z'))$$

does not evaluate $(\lambda y. y) (\lambda z. z)$ at all because it is not used in the evaluation:

$$\begin{aligned} & \frac{(\lambda t. \lambda f. f) ((\lambda y. y) (\lambda z. z)) ((\lambda y'. y') (\lambda z'. z'))}{(\lambda f. f) ((\lambda y'. y') (\lambda z'. z'))} \\ \mapsto & \dots \end{aligned}$$

The call-by-value strategy evaluates $(\lambda y. y) (\lambda z. z)$, but the result $\lambda z. z$ is ignored in the next reduction:

$$\begin{aligned} & (\lambda t. \lambda f. f) ((\lambda y. y) (\lambda z. z)) ((\lambda y'. y') (\lambda z'. z')) \\ \mapsto & \frac{(\lambda t. \lambda f. f) (\lambda z. z) ((\lambda y'. y') (\lambda z'. z'))}{(\lambda f. f) ((\lambda y'. y') (\lambda z'. z'))} \\ \mapsto & \dots \end{aligned}$$

The call-by-name strategy is adopted by the functional language Haskell. Haskell is called a *lazy* or *non-strict* functional language because it evaluates arguments to functions only if necessary (i.e., “lazily”). The actual implementation of Haskell uses another reduction strategy called *call-by-need*, which is semantically equivalent to the call-by-name strategy but never evaluates the same expression more than once. The call-by-value strategy is adopted by SML which is called an *eager* or *strict* functional language because it always evaluates arguments to functions regardless of whether they are actually used in function bodies or not (i.e., “eagerly”).

We say that an expression is in *normal form* if no reduction rule is applicable. Clearly every value (which is a λ -abstraction in the case of the λ -calculus) is in normal form. There are, however, expressions in normal form that are not values. For example, $x \lambda y. y$ is in normal form because x cannot be further reduced, but it is not a value either. We say that such expression is *stuck* or its reduction gets *stuck*. A stuck expression may be thought of as an ill-formed program, and ideally should not arise during an evaluation. Chapter ?? presents an extension of the λ -calculus which statically (i.e., at compile time) guarantees that a program satisfying a certain criterion never gets stuck.

1.3 Substitution

This section presents a definition of substitution $[e'/x]e$ to complete the operational semantics of the λ -calculus. While an informal interpretation of $[e'/x]e$ is obvious, its formal definition is a lot trickier than it appears.

First we need the notion of *free* variable which is the opposite of the notion of bound variable and plays a key role in the definition of substitution. A free variable is a variable that is not bound in any enclosing λ -abstraction. For example, y in $\lambda x. y$ is a free variable because no λ -abstraction of the form $\lambda y. e$ encloses its occurrence. To formalize the notion of free variable, we introduce a mapping $FV(e)$ to mean the set of free variables in e :

$$\begin{aligned} FV(x) &= \{x\} \\ FV(\lambda x. e) &= FV(e) - \{x\} \\ FV(e_1 e_2) &= FV(e_1) \cup FV(e_2) \end{aligned}$$

Since a variable is either free or bound, a variable x in e such that $x \notin FV(e)$ must be bound in some λ -abstraction. We say that an expression e is *closed* if it contains no free variables, i.e., $FV(e) = \emptyset$. Here are a few examples:

$$\begin{aligned} FV(\lambda x. x) &= \{\} \\ FV(x y) &= \{x, y\} \\ FV(\lambda x. x y) &= \{y\} \\ FV(\lambda y. \lambda x. x y) &= \{\} \\ FV((\lambda x. x y) (\lambda x. x z)) &= \{y, z\} \end{aligned}$$

A substitution $[e/x]e'$ is defined inductively with the following cases:

$$\begin{aligned} [e/x]x &= e \\ [e/x]y &= y && \text{if } x \neq y \\ [e/x](e_1 e_2) &= [e/x]e_1 [e/x]e_2 \end{aligned}$$

In order to give the definition of the remaining case $[e/x]\lambda y. e'$, we need to understand two properties of variables. The first property is that the name of a bound variable does not matter, which also conforms to our intuition. For example, an identity function $\lambda x. x$ inside an expression e may be rewritten as $\lambda y. y$ for an arbitrary variable y without changing the intended meaning of e , since both $\lambda x. x$ and $\lambda y. y$ denote an identity function. Another example is to rewrite $\lambda x. \lambda y. x y$ as $\lambda y. \lambda x. y x$ where both expressions denote the same function that applies the first argument to the second argument.

Formally we use a judgment $e \equiv_\alpha e'$ to mean that e can be rewritten as e' by renaming bound variables in e and vice versa. Here are examples of $e \equiv_\alpha e'$:

$$\begin{aligned} \lambda x. x &\equiv_\alpha \lambda y. y \\ \lambda x. \lambda y. x y &\equiv_\alpha \lambda z. \lambda y. z y \\ \lambda x. \lambda y. x y &\equiv_\alpha \lambda x. \lambda z. x z \\ \lambda x. \lambda y. x y &\equiv_\alpha \lambda y. \lambda x. y x \end{aligned}$$

By a historical accident, \equiv_α is called the α -equivalence relation, or we say that an α -conversion of e into e' rewrites e as e' by renaming bound variables in e . It turns out that a definition of $e \equiv_\alpha e'$ is also tricky to develop, which is given at the end of the present section.

The first property justifies the following case of substitution:

$$[e'/x]\lambda x. e = \lambda x. e$$

Intuitively, if we rewrite $\lambda x. e$ as another λ -abstraction of the form $\lambda y. e''$ where y is a fresh variable such that $x \neq y$, the substitution $[e'/x]$ is effectively ignored because x is found nowhere in $\lambda y. e''$. Here is a simple example with $e = x$:

$$\begin{aligned} [e'/x]\lambda x. x &\equiv_\alpha [e'/x]\lambda y. y \\ &= \lambda y. y \\ &\equiv_\alpha \lambda x. x \end{aligned}$$

A generalization of the case is that $[e'/x]$ has no effect on e if x is not a free variable in e :

$$[e'/x]e = e \quad \text{if } x \notin FV(e)$$

That is, we want to apply $[e'/x]$ to e only if x is a free variable in e .

The second property is that a free variable x in an expression e never turns into a bound variable; when explicitly replaced by another expression e' , as in $[e'/x]e$, it simply disappears. To better understand the second property, let us consider a naive definition of $[e'/x]\lambda y. e$ where y may or may not be a free variable in e' :

$$[e'/x]\lambda y. e = \lambda y. [e'/x]e \quad \text{if } x \neq y$$

Now, if y is a free variable in e' , it automatically becomes a bound variable in $\lambda y. [e'/x]e$, which is not acceptable. Here is an example showing such an anomaly:

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

Before the substitution, $\lambda y. x$ is a λ -abstraction that ignores its argument and returns x , but after the substitution, it turns into an identity function! What happens in the example is that a free variable y to be substituted for x is supposed to remain free after the substitution, but is accidentally *captured* by the λ -abstraction $\lambda y. x$ and becomes a bound variable. Such a phenomenon is called a *variable capture* which destroys the intuition that a free variable remains free unless it is replaced by another expression. This observation is generalized in the following definition of $[e'/x]\lambda y. e$ which is called a *capture-avoiding substitution*:

$$[e'/x]\lambda y. e = \lambda y. [e'/x]e \quad \text{if } x \neq y, y \notin FV(e')$$

If a variable capture occurs because $y \in FV(e')$, we rename y to another variable that is not free in e . For example, $(\lambda x. \lambda y. x y) y$ can be safely reduced after renaming the bound variable y to a fresh variable z :

$$(\lambda x. \lambda y. x y) y \mapsto [y/x] \lambda y. x y \equiv_{\alpha} [y/x] \lambda z. x z = \lambda z. y z$$

In the literature, the unqualified term “substitution” universally means a capture-avoiding substitution which renames bound variables as necessary.

Now we give a definition of the judgment $e \equiv_{\alpha} e'$. We need the notion of *variable swapping* $[x \leftrightarrow y]e$ which is obtained by replacing *all* occurrences of x in e by y and *all* occurrences of y in e by x . We emphasize that “all” occurrences include even those next to λ in λ -abstractions, which makes it straightforward to implement $[x \leftrightarrow y]e$. Here is an example:

$$[x \leftrightarrow y] \lambda x. \lambda y. x y = \lambda y. [x \leftrightarrow y] \lambda y. x y = \lambda y. \lambda x. [x \leftrightarrow y](x y) = \lambda y. \lambda x. y x$$

The definition of $e \equiv_{\alpha} e'$ is given inductively by the following inference rules:

$$\boxed{\begin{array}{c} \frac{}{x \equiv_{\alpha} x} \text{Var}_{\alpha} \quad \frac{e_1 \equiv_{\alpha} e'_1 \quad e_2 \equiv_{\alpha} e'_2}{e_1 e_2 \equiv_{\alpha} e'_1 e'_2} \text{App}_{\alpha} \\[10pt] \frac{e \equiv_{\alpha} e'}{\lambda x. e \equiv_{\alpha} \lambda x. e'} \text{Lam}_{\alpha} \quad \frac{x \neq y \quad y \notin FV(e) \quad [x \leftrightarrow y]e \equiv_{\alpha} e'}{\lambda x. e \equiv_{\alpha} \lambda y. e'} \text{Lam}'_{\alpha} \end{array}}$$

The rule Lam_{α} says that to compare $\lambda x. e$ and $\lambda x. e'$ which bind the same variable, we compare their bodies e and e' . To compare two λ -abstractions binding different variables, we use the rule Lam'_{α} .

To see why the rule Lam'_{α} works, we need to understand the implication of the premise $y \notin FV(e)$. Since $y \notin FV(e)$ implies $y \notin FV(\lambda x. e)$ and we have $x \notin FV(\lambda x. e)$, an outside observer would notice no difference even if the two variables x and y were literally swapped in $\lambda x. e$. In other words, $\lambda x. e$ and $[x \leftrightarrow y] \lambda x. e$ are effectively the same from the point of view of an outside observer. Since $[x \leftrightarrow y] \lambda x. e = \lambda y. [x \leftrightarrow y]e$, we compare $[x \leftrightarrow y]e$ with e' , which is precisely the third premise in the rule Lam'_{α} . As an example, here is a proof of $\lambda x. \lambda y. x y \equiv_{\alpha} \lambda y. \lambda x. y x$:

$$\frac{\frac{\frac{y \equiv_{\alpha} y}{y \equiv_{\alpha} y} \text{Var}_{\alpha} \quad \frac{x \equiv_{\alpha} x}{x \equiv_{\alpha} x} \text{Var}_{\alpha}}{y x \equiv_{\alpha} y x} \text{App}_{\alpha}}{\lambda x. y x \equiv_{\alpha} \lambda x. y x} \text{Lam}_{\alpha} \quad \frac{x \neq y \quad y \notin FV(\lambda y. x y)}{\lambda x. \lambda y. x y \equiv_{\alpha} \lambda y. \lambda x. y x} \text{Lam}'_{\alpha}$$

Exercise 1.1. Can we prove $\lambda x. e \equiv_{\alpha} \lambda y. e'$ when $x \neq y$ and $y \in FV(e)$?

Exercise 1.2. Suppose $x \notin FV(e)$ and $y \notin FV(e)$. Prove $e \equiv_{\alpha} [x \leftrightarrow y]e$.

Finally we give a complete definition of substitution:

$$\boxed{\begin{array}{ll} [e/x]x &= e \\ [e/x]y &= y & \text{if } x \neq y \\ [e/x](e_1 e_2) &= [e/x]e_1 [e/x]e_2 \\ [e'/x]\lambda x. e &= \lambda x. e \\ [e'/x]\lambda y. e &= \lambda y. [e'/x]e & \text{if } x \neq y, y \notin FV(e') \\ [e'/x]\lambda y. e &= \lambda z. [e'/x][y \leftrightarrow z]e & \text{if } x \neq y, y \in FV(e') \\ & & \text{where } z \neq y, z \notin FV(e), z \neq x, z \notin FV(e') \end{array}}$$

The last equation implies that if y is a free variable in e' , we choose another variable z satisfying the *where* clause and rewrite $\lambda y. e$ as $\lambda z. [y \leftrightarrow z]e$ by α -conversion:

$$\frac{y \neq z \quad z \notin FV(e) \quad [y \leftrightarrow z]e \equiv_{\alpha} [y \leftrightarrow z]e}{\lambda y. e \equiv_{\alpha} \lambda z. [y \leftrightarrow z]e} \text{Lam}'_{\alpha}$$

Then $z \notin FV(e')$ allows us to rewrite $[e'/x]\lambda z. [y \leftrightarrow z]e$ as $\lambda z. [e'/x][y \leftrightarrow z]e$. In a typical implementation, we obtain such a variable z just by generating a fresh variable. In such a case, replacing z by y never occurs and thus the last equation can be written as follows:

$$[e'/x]\lambda y. e = \lambda z. [e'/x][z/y]e$$

The new equation is less efficient, however. Consider $e = x y (\lambda y. x y)$ for example. $[y \leftrightarrow z]e$ gives $x z (\lambda z. x z)$ and $[e'/x][y \leftrightarrow z]e$ encounters no variable capture:

$$[e'/x][y \leftrightarrow z]e = [e'/x](x z (\lambda z. x z)) = e' z (\lambda z. e' z)$$

In contrast, $[z/y]e$ gives $x z (\lambda y. x y)$ and $[e'/x][z/y]e$ again encounters a variable capture in $[e'/x]\lambda y. x y$:

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

So we have to generate another fresh variable!

Exercise 1.3. What is the result of α -converting each expression in the left where a fresh variable to be generated in the conversion is provided in the right? Which expression is impossible to α -convert?

$$\begin{array}{ll} \lambda x. \lambda x'. x x' & \equiv_{\alpha} \lambda x'. \underline{\hspace{2cm}} \\ \lambda x. \lambda x'. x x' x'' & \equiv_{\alpha} \lambda x'. \underline{\hspace{2cm}} \\ \lambda x. \lambda x'. x x' x'' & \equiv_{\alpha} \lambda x''. \underline{\hspace{2cm}} \end{array}$$

1.4 Programming in the λ -calculus

In order to develop the λ -calculus to a full-fledged functional language, we need to show how to encode common datatypes such as boolean values, integers, and lists in the λ -calculus. Since all values in the λ -calculus are λ -abstractions, all such datatypes are also encoded with λ -abstractions. Once we show how to encode specific datatypes, we may use them as if they were built-in datatypes.

1.4.1 Church booleans

The inherent capability of a boolean value is to choose one of two different options. For example, a boolean truth chooses the first of two different options, as in an SML expression `if true then e_1 else e_2` . Thus boolean values in the λ -calculus, called Church booleans, are written as follows:

$$\begin{array}{ll} \text{tt} & = \lambda t. \lambda f. t \\ \text{ff} & = \lambda t. \lambda f. f \end{array}$$

Then a conditional construct `if e then e_1 else e_2` is defined as follows:

$$\text{if } e \text{ then } e_1 \text{ else } e_2 = e e_1 e_2$$

Here are examples of reducing conditional constructs under the call-by-name strategy:

$$\begin{array}{l} \text{if tt then } e_1 \text{ else } e_2 = \text{tt } e_1 e_2 = (\lambda t. \lambda f. t) e_1 e_2 \mapsto (\lambda f. e_1) e_2 \mapsto e_1 \\ \text{if ff then } e_1 \text{ else } e_2 = \text{ff } e_1 e_2 = (\lambda t. \lambda f. f) e_1 e_2 \mapsto (\lambda f. f) e_2 \mapsto e_2 \end{array}$$

Logical operators on boolean values are defined as follows:

$$\begin{array}{ll} \text{and} & = \lambda x. \lambda y. x y \text{ ff} \\ \text{or} & = \lambda x. \lambda y. x \text{ tt } y \\ \text{not} & = \lambda x. x \text{ ff } \text{tt} \end{array}$$

As an example, here are sequences of reductions of `and $e_1 e_2$` when $e_1 \mapsto^* \text{tt}$ and $e_1 \mapsto^* \text{ff}$, respectively, under the call-by-name strategy:

$$\begin{array}{ll} \text{and } e_1 e_2 & \text{and } e_1 e_2 \\ \mapsto^* e_1 e_2 \text{ ff} & \mapsto^* e_1 e_2 \text{ ff} \\ \mapsto^* \text{tt } e_2 \text{ ff} & \mapsto^* \text{ff } e_2 \text{ ff} \\ \mapsto^* e_2 & \mapsto^* \text{ff} \end{array}$$

The left sequence shows that when $e_1 \mapsto^* \text{tt}$ holds, and $e_1 e_2$ denotes the same truth value as e_2 . The right sequence shows that when $e_1 \mapsto^* \text{ff}$ holds, and $e_1 e_2$ evaluates to `ff` regardless of e_2 .

Exercise 1.4. Consider the conditional construct `if e then e1 else e2` defined as `e e1 e2` under the call-by-value strategy. How is it different from the conditional construct in SML?

Exercise 1.5. Define the logical operator `xor`. An easy way to define it is to use a conditional construct and the logical operator `not`.

1.4.2 Pairs

The inherent capability of a pair is to carry two unrelated values and to retrieve either value when requested. Thus, in order to represent a pair of e_1 and e_2 , we build a λ -abstraction which returns e_1 and e_2 when applied to `tt` and `ff`, respectively. Projection operators treat a pair as a λ -abstraction and applies it to either `tt` or `ff`.

$$\begin{aligned}\text{pair} &= \lambda x. \lambda y. \lambda b. b \ x \ y \\ \text{fst} &= \lambda p. p \ \text{tt} \\ \text{snd} &= \lambda p. p \ \text{ff}\end{aligned}$$

As an example, let us reduce `fst (pair e1 e2)` under the call-by-name strategy. Note that `pair e1 e2` evaluates to $\lambda b. b \ e_1 \ e_2$ which expects a boolean value for b in order to select either e_1 or e_2 . If `tt` is substituted for b , then $b \ e_1 \ e_2$ reduces to e_1 .

$$\begin{aligned}\text{fst} (\text{pair } e_1 \ e_2) &\mapsto (\text{pair } e_1 \ e_2) \ \text{tt} \\ &\mapsto^* (\lambda b. b \ e_1 \ e_2) \ \text{tt} \\ &\mapsto \text{tt } e_1 \ e_2 \\ &\mapsto^* e_1\end{aligned}$$

1.4.3 Church numerals

The inherent capability of a natural number n is to repeat a given process n times. In the case of the λ -calculus, n is encoded as a λ -abstraction \hat{n} , called a Church numeral, that takes a function f and returns $f^n = f \circ f \cdots \circ f$ (n times). Note that f^0 is an identity function $\lambda x. x$ because f is applied 0 times to its argument x , and that $\hat{1}$ itself is an identity function $\lambda f. f$.

$$\begin{aligned}\hat{0} &= \lambda f. f^0 = \lambda f. \lambda x. x \\ \hat{1} &= \lambda f. f^1 = \lambda f. \lambda x. f \ x \\ \hat{2} &= \lambda f. f^2 = \lambda f. \lambda x. f \ (f \ x) \\ \hat{3} &= \lambda f. f^3 = \lambda f. \lambda x. f \ (f \ (f \ x)) \\ &\vdots \\ \hat{n} &= \lambda f. f^n = \lambda f. \lambda x. f \ (f \ (f \ \cdots (f \ x) \cdots))\end{aligned}$$

If we read f as `S` and x as `0`, $\hat{n} \ f \ x$ returns the representation of the natural number n shown in Chapter ??.

Now let us define arithmetic operations on natural numbers. The addition operation `add \hat{m} \hat{n}` returns $\widehat{m+n}$ which is a λ -abstraction taking a function f and returning f^{m+n} . Since f^{m+n} may be written as $\lambda x. f^{m+n} \ x$, we develop `add` as follows; in order to differentiate natural numbers (e.g., n) from their encoded form (e.g., \hat{n}), we use \hat{m} and \hat{n} as variables:

$$\begin{aligned}\text{add} &= \lambda \hat{m}. \lambda \hat{n}. \lambda f. f^{m+n} \\ &= \lambda \hat{m}. \lambda \hat{n}. \lambda f. \lambda x. f^{m+n} \ x \\ &= \lambda \hat{m}. \lambda \hat{n}. \lambda f. \lambda x. f^m \ (f^n \ x) \\ &= \lambda \hat{m}. \lambda \hat{n}. \lambda f. \lambda x. \hat{m} \ f \ (\hat{n} \ f \ x)\end{aligned}$$

Note that f^m is obtained as $\hat{m} \ f$ (and similarly for f^n).

Exercise 1.6. Define the multiplication operation `mult \hat{m} \hat{n}` which returns $\widehat{m * n}$.

The multiplication operation can be defined in two ways. An easy way is to exploit the equation $m * n = m + m + \cdots + m$ (n times). That is, $m * n$ is obtained by adding m to zero exactly n times.

Since $\text{add } \hat{m}$ is conceptually a function adding m to its argument, we apply $\text{add } \hat{m}$ to $\hat{0}$ exactly n times to obtain $\widehat{m * n}$, or equivalently apply $(\text{add } \hat{m})^n$ to $\hat{0}$:

$$\begin{aligned}\text{mult} &= \lambda \hat{m}. \lambda \hat{n}. (\text{add } \hat{m})^n \hat{0} \\ &= \lambda \hat{m}. \lambda \hat{n}. \hat{n} (\text{add } \hat{m}) \hat{0}\end{aligned}$$

An alternative way (which may in fact be easier to figure out than the first solution) is to exploit the equation $f^{m*n} = (f^m)^n = (\hat{m} f)^n = \hat{n} (\hat{m} f)$:

$$\text{mult} = \lambda \hat{m}. \lambda \hat{n}. \lambda f. \hat{n} (\hat{m} f)$$

The subtraction operation is more difficult to define than the previous two operations. Suppose that we have a predecessor function pred computing the predecessor of a given natural number: $\text{pred } \hat{n}$ returns $\widehat{n-1}$ if $n > 0$ and $\hat{0}$ otherwise. To define the subtraction operation $\text{sub } \hat{m} \hat{n}$ which returns $\widehat{m-n}$ if $m > n$ and $\hat{0}$ otherwise, we apply pred to \hat{m} exactly n times:

$$\begin{aligned}\text{sub} &= \lambda \hat{m}. \lambda \hat{n}. \text{pred}^n \hat{m} \\ &= \lambda \hat{m}. \lambda \hat{n}. (\hat{n} \text{ pred}) \hat{m}\end{aligned}$$

Exercise 1.7. Define the predecessor function pred . Use an idea similar to the one used in a tail-recursive implementation of the Fibonacci function.

The predecessor function pred uses an auxiliary function next which takes pair $\hat{k} \hat{m}$, ignores \hat{k} , and returns pair $\widehat{m} \widehat{m+1}$:

$$\text{next} = \lambda p. \text{pair} (\text{snd } p) (\text{add} (\text{snd } p) \hat{1})$$

It can be shown that by applying next to pair $\hat{0} \hat{0}$ exactly n times, we obtain pair $\widehat{n-1} \hat{n}$ if $n > 0$ (under a certain reduction strategy):

$$\begin{aligned}\text{next}^0 (\text{pair } \hat{0} \hat{0}) &\mapsto^* \text{pair } \hat{0} \hat{0} \\ \text{next}^1 (\text{pair } \hat{0} \hat{0}) &\mapsto^* \text{pair } \hat{0} \hat{1} \\ \text{next}^2 (\text{pair } \hat{0} \hat{0}) &\mapsto^* \text{pair } \hat{1} \hat{2} \\ &\vdots \\ \text{next}^n (\text{pair } \hat{0} \hat{0}) &\mapsto^* \text{pair } \widehat{n-1} \hat{n}\end{aligned}$$

Since the predecessor of 0 is 0 anyway, the first component of $\text{next}^n (\text{pair } \hat{0} \hat{0})$ encodes the predecessor of n . Thus pred is defined as follows:

$$\begin{aligned}\text{pred} &= \lambda \hat{n}. \text{fst} (\text{next}^n (\text{pair } \hat{0} \hat{0})) \\ &= \lambda \hat{n}. \text{fst} (\hat{n} \text{ next} (\text{pair } \hat{0} \hat{0}))\end{aligned}$$

Exercise 1.8. Define a function $\text{isZero} = \lambda \hat{n}. \dots$ which tests if a given Church numeral is $\hat{0}$. Use it to define another function $\text{eq} = \lambda \hat{m}. \lambda \hat{n}. \dots$ which tests if two given Church numerals are equal.

1.5 Fixed point combinator

Since the λ -calculus is equally powerful as Turing machines, every Turing machine can be simulated by a certain expression in the λ -calculus. In particular, there are expressions in the λ -calculus that correspond to Turing machines that do not terminate and Turing machines that compute *recursive functions*.

It is relatively easy to find an expression whose reduction does not terminate. Suppose that we wish to find an expression ω such that $\omega \mapsto \omega$. Since it reduces to the same expression, its reduction never terminates. We rewrite ω as $(\lambda x. e) e'$ so that the β -reduction can be applied to the whole expression ω . Then we have

$$\omega = (\lambda x. e) e' \mapsto [e'/x]e = \omega.$$

Now $\text{omega} = [e'/x]e = (\lambda x. e) e'$ suggests $e = e'' x$ for some expression e'' such that $[e'/x]e'' = \lambda x. e$ (and $[e'/x]x = e'$):

$$\begin{aligned} \text{omega} &= [e'/x]e \\ &= [e'/x](e'' x) && \text{from } e = e'' x \\ &= [e'/x]e'' [e'/x]x \\ &= [e'/x]e'' e' \\ &= (\lambda x. e) e' && \text{from } [e'/x]e'' = \lambda x. e \end{aligned}$$

From $e = e'' x$ and $[e'/x]e'' = \lambda x. e$, we obtain $[e'/x]e'' = \lambda x. e'' x$. By letting $e'' = x$ in $[e'/x]e'' = \lambda x. e'' x$, we obtain $e' = \lambda x. x x$. Then omega can be defined as follows:

$$\begin{aligned} \text{omega} &= (\lambda x. e) e' \\ &= (\lambda x. e'' x) e' && \text{from } e = e'' x \\ &= (\lambda x. x x) e' && \text{from } e'' = x \\ &= (\lambda x. x x) (\lambda x. x x) && \text{from } e' = \lambda x. x x \end{aligned}$$

Now it can be shown that the reduction of omega defined as above never terminates.

Then how do we write recursive functions in the λ -calculus? We begin by assuming a *recursive function construct* $\text{fun } f \ x. e$ which defines a recursive function f whose argument is x and whose body is e . Note that the body e may contain references to f . Our goal is to show that $\text{fun } f \ x. e$ is *syntactic sugar* (which dissolves in the λ -calculus) in the sense that it can be rewritten as an existing expression in the λ -calculus and thus its addition does not increase the expressive power of the λ -calculus.

As a working example, we use a factorial function fac :

$$\text{fac} = \text{fun } f \ n. \text{if eq } n \ \hat{0} \text{ then } \hat{1} \text{ else mult } n \ (f \ (\text{pred } n))$$

Semantically f in the body refers to the very function fac being defined. First we mechanically derive a λ -abstraction $\text{FAC} = \lambda f. \lambda n. e$ from $\text{fac} = \text{fun } f \ n. e$:

$$\text{FAC} = \lambda f. \lambda n. \text{if eq } n \ \hat{0} \text{ then } \hat{1} \text{ else mult } n \ (f \ (\text{pred } n))$$

Note that FAC has totally different characteristics than fac : while fac takes a natural number n to return another natural number, FAC takes a function f to return another function. (If fac and FAC were allowed to have types, fac would have type $\text{nat} \rightarrow \text{nat}$ whereas FAC would have type $(\text{nat} \rightarrow \text{nat}) \rightarrow (\text{nat} \rightarrow \text{nat})$.)

The key idea behind constructing FAC is that given a *partial* implementation f of the factorial function, $\text{FAC } f$ returns an *improved* implementation of the factorial function. Suppose that f correctly computes the factorial of any natural number up to n . Then $\text{FAC } f$ correctly computes the factorial of any natural number up to $n + 1$, which is an improvement over f . Note also that $\text{FAC } f$ correctly computes the factorial of 0 regardless of f . In particular, even when given a least informative function $f = \lambda n. \text{omega}$ (which does nothing because it never returns), $\text{FAC } f$ correctly computes the factorial of 0. Thus we can imagine an infinite chain of functions $\{\text{fac}_0, \text{fac}_1, \dots, \text{fac}_i, \dots\}$ which begins with $\text{fac}_0 = \text{FAC } \lambda n. \text{omega}$ and repeatedly applies the equation $\text{fac}_{i+1} = \text{FAC } \text{fac}_i$:

$$\begin{aligned} \text{fac}_0 &= \text{FAC } \lambda n. \text{omega} \\ \text{fac}_1 &= \text{FAC } \text{fac}_0 = \text{FAC}^2 \lambda n. \text{omega} \\ \text{fac}_2 &= \text{FAC } \text{fac}_1 = \text{FAC}^3 \lambda n. \text{omega} \\ &\vdots \\ \text{fac}_i &= \text{FAC } \text{fac}_{i-1} = \text{FAC}^{i+1} \lambda n. \text{omega} \\ &\vdots \end{aligned}$$

Note that fac_i correctly computes the factorial of any natural number up to i . Then, if ω denotes an infinite natural number (greater than any natural number), we may take fac_ω as a correct implementation of the factorial function fac , i.e., $\text{fac} = \text{fac}_\omega$.

Another important observation is that given a correct implementation fac of the factorial function, $\text{FAC } \text{fac}$ returns another correct implementation of the factorial function. That is, if fac is a correct implementation of the factorial function,

$$\lambda n. \text{if eq } n \ \hat{0} \text{ then } \hat{1} \text{ else mult } n \ (\text{fac } (\text{pred } n))$$

is also a correct implementation of the factorial function. Since the two functions are essentially identical in that both return the same result for any argument, we may let $\text{fac} = \text{FAC fac}$. If we substitute fac_ω for fac in the equation, we obtain $\text{fac}_\omega = \text{fac}_{\omega+1}$ which also makes sense because $\omega \leq \omega + 1$ by the definition of $+$ and $\omega + 1 \leq \omega$ by the definition of ω (which is greater than any natural number including $\omega + 1$).

Now it seems that FAC contains all necessary information to derive $\text{fac} = \text{fac}_\omega = \text{FAC fac}$, but exactly how? It turns out that fac is obtained by applying the *fixed point combinator* fix to FAC , i.e., $\text{fac} = \text{fix FAC}$, where fix is defined as follows:

$$\text{fix} = \lambda F. (\lambda f. F (\lambda x. f f x)) (\lambda f. F (\lambda x. f f x))$$

Here we assume the call-by-value strategy; for the call-by-name strategy, we simplify $\lambda x. f f x$ into $f f$ and use the following fixed point combinator fix_{CBN} :

$$\text{fix}_{\text{CBN}} = \lambda F. (\lambda f. F (f f)) (\lambda f. F (f f))$$

To understand how the fixed point combinator fix works, we need to learn the concept of *fixed point*.² A fixed point of a function f is a value v such that $v = f(v)$. For example, the fixed point of a function $f(x) = 2 - x$ is 1 because $1 = f(1)$. As its name suggests, fix takes a function F (which itself transforms a function f into another function f') and returns its fixed point. That is, $\text{fix } F$ is a fixed point of F :

$$\text{fix } F = F (\text{fix } F)$$

Informally the left expression transforms into the right expression via the following steps; we use a symbol \approx to emphasize “informally” because the transformation is not completely justified by the β -reduction alone:

$$\begin{aligned} \text{fix } F &\mapsto g g && \text{where } g = \lambda f. F (\lambda x. f f x) \\ &= (\lambda f. F (\lambda x. f f x)) g \\ &\mapsto F (\lambda x. g g x) \\ &\approx F (g g) && \text{because } \lambda x. g g x \approx g g \\ &\approx F (\text{fix } F) && \text{because } \text{fix } F \mapsto g g \end{aligned}$$

Now we can explain why fix FAC gives an implementation of the factorial function. By the nature of the fixed point combinator fix , we have

$$\text{fix FAC} = \text{FAC} (\text{fix FAC}).$$

That is, fix FAC returns a function f satisfying $f = \text{FAC} f$, which is precisely the property that fac needs to satisfy! Therefore we take fix FAC as an equivalent of fac .³

An alternative way to explain the behavior of fix FAC is as follows. Suppose that we wish to compute $\text{fac } n$ for an arbitrary natural number n . Since fix FAC is a fixed point of FAC , we have the following equation:

$$\begin{aligned} \text{fix FAC} &= \text{FAC} (\text{fix FAC}) \\ &= \text{FAC} (\text{FAC} (\text{fix FAC})) = \text{FAC}^2 (\text{fix FAC}) \\ &= \text{FAC}^2 (\text{FAC} (\text{fix FAC})) = \text{FAC}^3 (\text{fix FAC}) \\ &\vdots \\ &= \text{FAC}^n (\text{FAC} (\text{fix FAC})) = \text{FAC}^{n+1} (\text{fix FAC}) \end{aligned}$$

The key observation is that $\text{FAC}^{n+1} (\text{fix FAC})$ correctly computes the factorial of any natural number up to n *regardless of what* fix FAC *does* (see Page 11). Since we have $\text{fix FAC} = \text{FAC}^{n+1} (\text{fix FAC})$, it follows that fix FAC correctly computes the factorial of an *arbitrary* natural number. That is, fix FAC does precisely what fac does.

In summary, in order to encode a recursive function $\text{fun } f x. e$ in the λ -calculus, we first derive a λ -abstraction $F = \lambda f. \lambda x. e$. Then $\text{fix } F$ *automagically* returns a function that exhibits the same behavior as $\text{fun } f x. e$ does.

Exercise 1.9. Under the call-by-value strategy, fac , or equivalently fix FAC , never terminates when applied to any natural number! Why? (Hint: Exercise 1.4)

²Never use the word *fixpoint*! Dana Scott, who coined the word *fixed point*, says that *fixpoint* is wrong!

³The fixed point combinator fix actually yields what is called the *least* fixed point. That is, a function F may have many fixed points and fix returns the least one in the sense that the least one is the most informative one. The least fixed point is what we usually expect.

1.6 Deriving the fixed point combinator

This section explains how to derive the fixed point combinator. As its formal derivation is extremely intricate, we will illustrate the key idea with an example. Students may choose to skip this section if they wish.

Let us try to write a factorial function fac without using the fixed point combinator. Consider the following function $\text{fac}_{\text{wrong}}$:

$$\text{fac}_{\text{wrong}} = \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (f (\text{pred } n))$$

$\text{fac}_{\text{wrong}}$ is simply wrong because its body contains a reference to an unbound variable f . If, however, f points to a correct implementation fac of the factorial function, $\text{fac}_{\text{wrong}}$ would also be a correct implementation. Since there is no way to use a free variable f in reducing an expression, we have to introduce it in a λ -abstraction anyway:

$$\text{FAC} = \lambda f. \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (f (\text{pred } n))$$

FAC is definitely an improvement over $\text{fac}_{\text{wrong}}$, but it is not a function taking a natural number; rather it takes a function f to return another function which refines f . More importantly, there seems to be no way to make a recursive call with FAC because FAC calls only its argument f in its body and never makes a recursive call to itself.

Then how do we make a recursive call with FAC? The problem at hand is that the body of FAC, which needs to call fac , calls only its argument f . Our instinct, however, says that FAC contains all necessary information to derive fac (i.e., $\text{FAC} \approx \text{fac}$) because its body resembles a typical implementation of the factorial function. Thus we are led to try substituting FAC itself for f . That is, we make a call to FAC using FAC itself as an argument — what a crazy idea it is!:

$$\text{FAC FAC} = \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (\text{FAC} (\text{pred } n))$$

Unfortunately FAC FAC returns a function which does not make sense: in its body, a call to FAC is made with an argument $\text{pred } n$, but FAC expects not a natural number but a function. It is, however, easy to fix the problem: if FAC FAC returns a correct implementation of the factorial function, we only need to replace FAC in the body by FAC FAC. That is, what we want in the end is the following equation

$$\text{FAC FAC} = \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (\text{FAC FAC} (\text{pred } n))$$

where FAC FAC serves as a correct implementation of the factorial function.

Let us change the definition of FAC so that it satisfies the above equation. All we need to do is to replace a reference to f in its body by an application $f f$. Thus we obtain a new function Fac defined as follows:

$$\text{Fac} = \lambda f. \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (f f (\text{pred } n))$$

It is easy to see that Fac satisfies the following equation:

$$\text{Fac Fac} = \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (\text{Fac Fac} (\text{pred } n))$$

Since Fac Fac returns a correct implementation of the factorial function, we define fac as follows:

$$\text{fac} = \text{Fac Fac}$$

Now let us derive the fixed point combinator fix by rewriting fac in terms of fix (and FAC as it turns out). Consider the body of Fac:

$$\text{Fac} = \lambda f. \underline{\lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (f f (\text{pred } n))}$$

The underlined expression is almost the body of a typical factorial function except for the application $f f$. The following definition of Fac abstracts from the application $f f$ by replacing it by a reference to a single function g :

$$\begin{aligned} \text{Fac} &= \lambda f. (\lambda g. \lambda n. \text{if eq } n \hat{0} \text{ then } \hat{1} \text{ else mult } n (g (\text{pred } n))) (f f) \\ &= \lambda f. \underline{\text{FAC}} (f f) \end{aligned}$$

Then `fac` is rewritten as follows:

$$\begin{aligned}
\text{fac} &= \text{Fac Fac} \\
&= (\lambda f. \text{FAC } (f f)) (\lambda f. \text{FAC } (f f)) \\
&= \lambda F. ((\lambda f. F (f f)) (\lambda f. F (f f))) \text{FAC} \\
&= \text{fix}_{\text{CBN}} \text{FAC}
\end{aligned}$$

In the case of the call-by-value strategy, `fixCBN FAC` always diverges. A quick fix is to rewrite `f f` as `λx. f f x` and we obtain `fix`:

$$\begin{aligned}
\text{fac} &= \text{Fac Fac} \\
&= \lambda F. ((\lambda f. F (f f)) (\lambda f. F (f f))) \text{FAC} \\
&= \lambda F. ((\lambda f. F (\lambda x. f f x)) (\lambda f. F (\lambda x. f f x))) \text{FAC} \\
&= \text{fix FAC}
\end{aligned}$$

This is how to derive the fixed point combinator `fix`!

1.7 De Bruijn indexes

As λ -abstractions are intended to denote mathematical functions with formal arguments, variable names may seem to be an integral part of the syntax for the λ -calculus. For example, it seems inevitable to introduce a formal argument, say x , when defining an identity function. On the other hand, a specific choice of a variable name does not affect the meaning of a λ -abstraction. For example, $\lambda x. x$ and $\lambda y. y$ both denote the same identity function even though they bind different variable names as formal arguments. In general, α -conversion enables us to rewrite any λ -abstraction into another λ -abstraction with a different name for the bound variable. This observation suggests that there may be a way to represent variables in the λ -calculus without specific names. An example of such a nameless representation of variables is *de Bruijn indexes*.

The basic idea behind de Bruijn indexes is to represent each variable by an integer value, called a de Bruijn index, instead of a name. (De Bruijn indexes can be negative, but we consider non-negative indexes only.) Roughly speaking, a de Bruijn index counts the number of λ -binders, such as λx , λy , and λz , lying between a given variable and its corresponding (unique) λ -binder. For example, x in the body of $\lambda x. x$ is assigned a de Bruijn index 0 because there is no intervening λ -binder between x and λx . In contrast, x in the body of $\lambda x. \lambda y. x y$ is assigned a de Bruijn index 1 because there lies an intervening λ -binder λy between x and λx . Thus a de Bruijn index for a variable specifies the relative position of its corresponding λ -binder. This, in turn, implies that the same variable can be assigned different de Bruijn indexes depending on its position. For example, in $\lambda x. x (\lambda y. x y)$, the first occurrence of x is assigned 0 whereas the second occurrence is assigned 1 because of the λ -binder λy .

Since all variables are now represented by integer values, there is no need to explicitly introduce variables in λ -abstractions. In fact, it is impossible because the same variable can be assigned different de Bruijn indexes. Thus, expressions with de Bruijn indexes, or *de Bruijn expressions*, are inductively defined as follows:

$$\begin{array}{ll}
\text{de Bruijn expression} & M ::= n \mid \lambda. M \mid M M \\
\text{de Bruijn index} & n ::= 0 \mid 1 \mid 2 \mid \dots
\end{array}$$

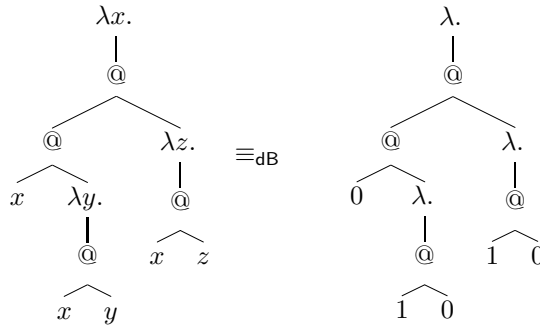
For de Bruijn expressions, we use metavariables M and N ; for de Bruijn indexes, we use metavariables n , m , and i .

We write $e \equiv_{\text{dB}} M$ to mean that an ordinary expression e is converted to a de Bruijn expression M . Sometimes it suffices to literally count the number of λ -binders lying between each variable and its corresponding λ -binder, as in all the examples given above:

$$\begin{aligned}
\lambda x. x &\equiv_{\text{dB}} \lambda. 0 \\
\lambda x. \lambda y. x y &\equiv_{\text{dB}} \lambda. \lambda. 1 0 \\
\lambda x. x (\lambda y. x y) &\equiv_{\text{dB}} \lambda. 0 (\lambda. 1 0)
\end{aligned}$$

In general, however, converting an ordinary expression e into a de Bruijn expression requires us to interpret e as a tree-like structure rather than a linear structure. As an example, consider $\lambda x. (x (\lambda y. x y)) (\lambda z. x z)$. Literally counting the number of λ -binders results in a de Bruijn expression $\lambda. (0 (\lambda. 1 0)) (\lambda. 2 0)$, in

which the last occurrence of x is assigned a (wrong) de Bruijn index 2 because of λy and λz . Intuitively, however, the last occurrence of x must be assigned a de Bruijn index 1 because its corresponding λ -binder can be located irrespective of the λ -binder λy . Thus, a proper way to convert an expression e to a de Bruijn expression is to count the number of λ -binders found along the way from each variable to its corresponding λ -binder in the tree-like representation of e . For example, we have $\lambda x. (x (\lambda y. x y)) (\lambda z. x z) \equiv_{dB} \lambda. (0 (\lambda. 1 0)) (\lambda. 1 0)$ as illustrated below:



1.7.1 Substitution

In order to exploit de Bruijn indexes in implementing the operational semantics of the λ -calculus, we need a definition of substitution for de Bruijn expressions, from which a definition of β -reduction can be derived. We wish to define $\sigma_0(M, N)$ such that the following relationship holds:

$$\begin{array}{ccc} (\lambda x. e) e' & \mapsto & [e'/x]e \\ \equiv_{dB} & & \equiv_{dB} \\ (\lambda. M) N & \mapsto & \sigma_0(M, N) \end{array}$$

That is, applying $\lambda. M$ to N , or substituting N for the variable bound in $\lambda. M$, results in $\sigma_0(M, N)$. (The meaning of the subscript 0 in $\sigma_0(M, N)$ is explained later.)

Instead of beginning with a complete definition of $\sigma_0(M, N)$, let us refine it through a series of examples. Consider the following example in which the redex is underlined:

$$\begin{array}{ccc} \lambda x. \lambda y. \underline{(\lambda z. x y z)} (\lambda w. w) & \mapsto & \lambda x. \lambda y. x y (\lambda w. w) \\ \equiv_{dB} & & \equiv_{dB} \\ \lambda. \lambda. \underline{(\lambda. 2 1 0)} (\lambda. 0) & \mapsto & \lambda. \lambda. 1 0 (\lambda. 0) \end{array}$$

We observe that 0, which corresponds to z bound in the λ -abstraction $\lambda z. x y z$, is replaced by the argument $\lambda. 0$. The other indexes 1 and 2 are decremented by one because the λ -binder λz disappears. These two observations lead to the following partial definition of $\sigma_0(M, N)$:

$$\begin{array}{ll} \sigma_0(M_1 M_2, N) & = \sigma_0(M_1, N) \sigma_0(M_2, N) \\ \sigma_0(0, N) & = N \\ \sigma_0(m, N) & = m - 1 \quad \text{if } m > 0 \end{array}$$

To see how the remaining case $\sigma_0(\lambda. M, N)$ is defined, consider another example in which the redex is underlined:

$$\begin{array}{ccc} \lambda x. \lambda y. \underline{(\lambda z. (\lambda u. x y z u))} (\lambda w. w) & \mapsto & \lambda x. \lambda y. (\lambda u. x y (\lambda w. w) u) \\ \equiv_{dB} & & \equiv_{dB} \\ \lambda. \lambda. \underline{(\lambda. (\lambda. 3 2 1 0))} (\lambda. 0) & \mapsto & \lambda. \lambda. (\lambda. 2 1 (\lambda. 0) 0) \end{array}$$

We observe that unlike in the first example, 0 remains intact because it corresponds to u bound in $\lambda u. x y z u$, while 1 corresponds to z and is thus replaced by $\lambda. 0$. The reason why 1 is now replaced by $\lambda. 0$ is that in general, a de Bruijn index m outside $\lambda. M$ points to the same variable as $m + 1$ inside

$\lambda. M$, i.e., within M . This observation leads to an equation $\sigma_0(\lambda. M, N) = \lambda. \sigma_1(M, N)$ where $\sigma_1(M, N)$ is defined as follows:

$$\begin{aligned}\sigma_1(M_1 M_2, N) &= \sigma_1(M_1, N) \sigma_1(M_2, N) \\ \sigma_1(0, N) &= 0 \\ \sigma_1(1, N) &= N \\ \sigma_1(m, N) &= m - 1 \quad \text{if } m > 1\end{aligned}$$

In the two examples above, we see that the subscript n in $\sigma_n(M, N)$ serves as a “boundary” index: m remains intact if $m < n$, m is replaced by N if $m = n$, and m is decremented by one if $m > n$. Alternatively n in $\sigma_n(M, N)$ may be read as the number of λ -binders enclosing M as illustrated below:

$$\sigma_0(\underbrace{\lambda. \lambda. \dots \lambda.}_n M, N) = \underbrace{\lambda. \lambda. \dots \lambda.}_n \sigma_n(M, N)$$

The following definition of $\sigma_n(M, N)$ uses n as a boundary index and also generalizes the relationship between $\sigma_0(\lambda. M, N)$ and $\lambda. \sigma_1(M, N)$:

$$\begin{aligned}\sigma_n(M_1 M_2, N) &= \sigma_n(M_1, N) \sigma_n(M_2, N) \\ \sigma_n(\lambda. M, N) &= \lambda. \sigma_{n+1}(M, N) \\ \sigma_n(m, N) &= m \quad \text{if } m < n \\ \sigma_n(n, N) &= N \\ \sigma_n(m, N) &= m - 1 \quad \text{if } m > n\end{aligned}$$

The following example combines the two examples given above:

$$\begin{aligned}\lambda x. \lambda y. \underbrace{(\lambda z. (\lambda u. x y z u) (x y z)) (\lambda w. w)}_{\equiv_{dB}} &\mapsto \lambda x. \lambda y. (\lambda u. x y (\lambda w. w) u) (x y (\lambda w. w)) \\ \lambda. \lambda. \underbrace{(\lambda. (\lambda. 3 \ 2 \ 1 \ 0) (2 \ 1 \ 0)) (\lambda. 0)}_{\equiv_{dB}} &\mapsto \lambda. \lambda. (\lambda. 2 \ 1 (\lambda. 0) 0) (1 \ 0 (\lambda. 0))\end{aligned}$$

The use of de Bruijn indexes obviates the need for α -conversion because variable names never clash. Put simply, there is no need to rename bound variables to avoid variable captures because variables have no names anyway.

1.7.2 Shifting

Although the previous definition of $\sigma_n(M, N)$ is guaranteed to work if N is closed, it may not work if N represents an expression with free variables. To be specific, the equation $\sigma_n(n, N) = N$ ceases to hold if $n > 0$ and N represents an expression with free variables. Consider the following example in which the redex is underlined:

$$\begin{aligned}\lambda x. \lambda y. \underbrace{(\lambda z. (\lambda u. z) z) (\lambda w. x y w)}_{\equiv_{dB}} &\mapsto \lambda x. \lambda y. (\lambda u. \lambda w. x y w) (\lambda w. x y w) \\ \lambda. \lambda. \underbrace{(\lambda. (\lambda. 1) 0) (\lambda. 2 \ 1 \ 0)}_{\equiv_{dB}} &\mapsto \lambda. \lambda. (\lambda. \lambda. 3 \ 2 \ 0) (\lambda. 2 \ 1 \ 0)\end{aligned}$$

The previous definition of $\sigma_n(M, N)$ yields a wrong result because $\sigma_1(1, \lambda. 2 \ 1 \ 0)$ yields $\lambda. 2 \ 1 \ 0$ instead of $\lambda. 3 \ 2 \ 0$:

$$\begin{aligned}\underbrace{(\lambda. (\lambda. 1) 0) (\lambda. 2 \ 1 \ 0)} &\mapsto \sigma_0((\lambda. 1) 0, \lambda. 2 \ 1 \ 0) \\ &= \sigma_0(\lambda. 1, \lambda. 2 \ 1 \ 0) \sigma_0(0, \lambda. 2 \ 1 \ 0) \\ &= (\lambda. \sigma_1(1, \lambda. 2 \ 1 \ 0)) \sigma_0(0, \lambda. 2 \ 1 \ 0) \\ &= (\lambda. \lambda. 2 \ 1 \ 0) (\lambda. 2 \ 1 \ 0) \\ &\neq (\lambda. \lambda. 3 \ 2 \ 0) (\lambda. 2 \ 1 \ 0)\end{aligned}$$

To see why $\sigma_n(n, N) = N$ fails to hold in general, recall that the subscript n in $\sigma_n(n, N)$ denotes the number of λ -binders enclosing the de Bruijn index n :

$$\sigma_0(\underbrace{\lambda. \lambda. \dots \lambda.}_n n, N) = \underbrace{\lambda. \lambda. \dots \lambda.}_n \sigma_n(n, N)$$

Therefore all de Bruijn indexes in N corresponding to free variables must be shifted by n after the substitution so that they correctly skip those n λ -binders enclosing the de Bruijn index n . For example, we have:

$$\sigma_0(\underbrace{\lambda. \lambda. \dots \lambda.}_n n, m) = \underbrace{\lambda. \lambda. \dots \lambda.}_n \sigma_n(n, m) = \underbrace{\lambda. \lambda. \dots \lambda.}_n m + n$$

(Here $m + n$ is a single de Bruijn index adding m and n , not a composite de Bruijn index expression consisting of m , n , and $+$.)

Let us write $\tau^n(N)$ for shifting by n all de Bruijn indexes in N corresponding to free variables. Now we use $\sigma_n(n, N) = \tau^n(N)$ instead of $\sigma_n(n, N) = N$. A partial definition of $\tau^n(N)$ is given as follows:

$$\begin{aligned} \tau^n(N_1 N_2) &= \tau^n(N_1) \tau^n(N_2) \\ \tau^n(m) &= m + n \end{aligned}$$

The remaining case $\tau^n(\lambda. N)$, however, cannot be defined inductively in terms of $\tau^n(N)$, for example, like $\tau^n(\lambda. N) = \lambda. \tau^n(N)$. The reason is that within N , not every de Bruijn index corresponds to a free variable: 0 finds its corresponding λ -binder in $\lambda. N$ and thus must remain intact.

This observation suggests that we need to maintain another “boundary” index (similar to the boundary index n in $\sigma_n(M, N)$) in order to decide whether a given de Bruijn index corresponds to a free variable or not. For example, if the boundary index for $\lambda. N$ starts at 0, it increments to 1 within N . Thus we are led to use a general form $\tau_i^n(N)$ for shifting by n all de Bruijn indexes in N where a de Bruijn index m in N corresponding to free variables such that $m < i$ does not count as a free variable. Formally $\tau_i^n(N)$ is defined as follows:

$$\begin{aligned} \tau_i^n(N_1 N_2) &= \tau_i^n(N_1) \tau_i^n(N_2) \\ \tau_i^n(\lambda. N) &= \lambda. \tau_{i+1}^n(N) \\ \tau_i^n(m) &= m + n && \text{if } m \geq i \\ \tau_i^n(m) &= m && \text{if } m < i \end{aligned}$$

Accordingly the complete definition of $\sigma_n(M, N)$ is given as follows:

$$\begin{aligned} \sigma_n(M_1 M_2, N) &= \sigma_n(M_1, N) \sigma_n(M_2, N) \\ \sigma_n(\lambda. M, N) &= \lambda. \sigma_{n+1}(M, N) \\ \sigma_n(m, N) &= m && \text{if } m < n \\ \sigma_n(n, N) &= \tau_0^n(N) \\ \sigma_n(m, N) &= m - 1 && \text{if } m > n \end{aligned}$$

Now $(\lambda. (\lambda. 1) 0) (\lambda. 2 1 0)$ from the earlier example reduces correctly:

$$\begin{aligned} \underline{(\lambda. (\lambda. 1) 0) (\lambda. 2 1 0)} &\mapsto \sigma_0((\lambda. 1) 0, \lambda. 2 1 0) \\ &= \sigma_0(\lambda. 1, \lambda. 2 1 0) \sigma_0(0, \lambda. 2 1 0) \\ &= (\lambda. \sigma_1(1, \lambda. 2 1 0)) \sigma_0(0, \lambda. 2 1 0) \\ &= (\lambda. \tau_0^1(\lambda. 2 1 0)) \tau_0^0(\lambda. 2 1 0) \\ &= (\lambda. \lambda. \tau_1^1(2 1 0)) \lambda. \tau_1^0(2 1 0) \\ &= (\lambda. \lambda. (2 + 1) (1 + 1) 0) (\lambda. (2 + 0) (1 + 0) 0) \\ &= (\lambda. \lambda. 3 2 0) (\lambda. 2 1 0) \end{aligned}$$

When converting an ordinary expression e with free variables x_0, x_1, \dots, x_n into a de Bruijn expression, we may convert $\lambda x_0. \lambda x_1. \dots \lambda x_n. e$ instead, which effectively assigns de Bruijn indexes $0, 1, \dots, n$ to x_0, x_1, \dots, x_n , respectively. Then we can think of reducing e as reducing $\lambda x_0. \lambda x_1. \dots \lambda x_n. e$ where the n λ -binders are all ignored. In this way, we can exploit de Bruijn indexes in reducing expressions with free variables (or global variables).

1.8 Exercises

Exercise 1.10. We wish to develop a weird reduction strategy for the λ -calculus:

- Given an application $e_1 e_2$, we first reduce e_2 .
- After reducing e_2 to a value, we reduce e_1 .
- When e_1 reduces to a λ -abstraction, we apply the β -reduction.

Give the rules for the reduction judgment $e \mapsto e'$ under the weird reduction strategy.

Exercise 1.11. In a reduction sequence judgment $e \mapsto^* e'$, we use \mapsto^* for the reflexive and transitive closure of \mapsto . That is, $e \mapsto^* e'$ holds if $e \mapsto e_1 \mapsto \dots \mapsto e_n = e'$ where $n \geq 0$. Formally we use the following inductive definition:

$$\frac{}{e \mapsto^* e} \text{ Refl} \quad \frac{e \mapsto e'' \quad e'' \mapsto^* e'}{e \mapsto^* e'} \text{ Trans}$$

We would expect that $e \mapsto^* e'$ and $e' \mapsto^* e''$ together imply $e \mapsto^* e''$ because we obtain a proof of $e \mapsto^* e''$ simply by concatenating $e \mapsto e_1 \mapsto \dots \mapsto e_n = e'$ and $e' \mapsto e'_1 \mapsto \dots \mapsto e'_m = e''$:

$$e \mapsto e_1 \mapsto \dots \mapsto e_n = e' \mapsto e'_1 \mapsto \dots \mapsto e'_m = e''$$

Give a proof of this transitivity property of \mapsto^* : if $e \mapsto^* e'$ and $e' \mapsto^* e''$, then $e \mapsto^* e''$. To which judgment of $e \mapsto^* e'$ and $e' \mapsto^* e''$ do we have to apply rule induction?

Exercise 1.12. Define a function `double` = $\lambda \hat{n}. \dots$ for doubling a given natural number encoded as a Church numeral. Specifically `double` \hat{n} returns $\widehat{2 * n}$.

Exercise 1.13. Define an operation `halve` for halving a given natural number. Specifically `halve` \hat{n} returns $\widehat{n/2}$:

- `halve` $\widehat{2 * k}$ returns \hat{k} .
- `halve` $(\widehat{2 * k + 1})$ returns \hat{k} .

You may use `pair`, `fst`, and `snd` without expanding them into their definitions. You may also use `zero` for a natural number zero and `succ` for finding the successor to a given natural number:

$$\begin{aligned} \text{zero} &= \hat{0} = \lambda f. \lambda x. x \\ \text{succ} &= \lambda \hat{n}. \lambda f. \lambda x. \hat{n} f (f x) \end{aligned}$$