

Chapter 1

Introduction

1.1 Outline

1. Definition of a simple, typed lambda calculus
2. Context and definition of *effects* in programming languages
3. Examples of common effects
 - (a) IO
 - (b) mutable data, state
 - (c) exception
 - (d) nondeterminism
 - (e) non-termination / non-totality
 - (f) reader, writer, output
 - (g) continuation
4. Explanation of how effects are handled in ML's style
 - (a) call-by-value at runtime
 - (b) examples of effects: IO, exceptions
 - (c) effects are untyped
 - (d) simple exception-handling system (type/catch)
5. Explain why this style is not purely functional
 - (a) exposes *side-effects*, which are implicit effects not accessible by the programming language

1.2 A Simply Typed Lambda-Calculus Chapter 1. Introduction

lambda-calculus is a formal language for expressing computation. In this context, a *language* is defined to be a set of expression that are generated by syntactical rules. The lambda-calculus comes in many different variants, and here we will consider a basic variant of the *simply-typed lambda-calculus* in order to considering computation formally. Call our language Language-A.

1.2.1 Definition of Language-A

TODO: include let bindings for Language-A?

Syntax for Language-A

In Language-A, there are two forms used for constructing well-formed expressions: *terms* and *types*. They are expressed by the following syntax:

metavariable	constructors	name
$\langle term \rangle$	$\langle term-name \rangle$ $\langle name \rangle : \langle type \rangle \Rightarrow \langle term \rangle$ $\langle term \rangle \langle term \rangle$ $\text{left } \langle term \rangle$ $\text{right } \langle term \rangle$ $\text{match } \langle term \rangle \{ \text{left } \langle name \rangle \Rightarrow \langle term \rangle \mid \text{right } \langle name \rangle \Rightarrow \langle term \rangle \}$ $(\langle term \rangle, \langle term \rangle)$ $\text{first } \langle term \rangle$ $\text{second } \langle term \rangle$	atom function application left right split product first second
$\langle type \rangle$	$\langle type-name \rangle$ $\langle type \rangle \rightarrow \langle type \rangle$ $\langle type \rangle + \langle type \rangle$ $\langle type \rangle \times \langle type \rangle$	atom arrow sum product

(1.1)

where $\langle term-name \rangle$ and $\langle type-name \rangle$ are metavariables that range over an infinite collection of distinct names, and $\langle name \rangle$ ranges over all names. A given expression is well-formed with respect to Language-A if it is constructable by a series of applications of these rules. Note that this syntax does not specify any concrete terms or types. Language-A's syntax is defined abstractly in order for the definition its typing and semantics as simple as possible. For the sake of intuition however, these are some examples of the likes of which these meta-variables stand in for:

- $\langle term-name \rangle$: 0, 1, 2, true, false, "hello world", •.
- $\langle type-name \rangle$: natural, integer, boolean, string, void, unit.
- $\langle name \rangle$: x, y, this-is-a-constant



Typing Rules for Language-A

Terms and types are related by a *typing judgement*, by which a term is states to have a unique type. The judgement that a has type α is written as $a:\alpha$. In order to build typed terms using the constructors presented in 1.1, the types of complex terms are inferred from their sub-terms using inference rules making use of judgement contexts (i.e. collections of judgements). A statement of the form $\Gamma \vdash a:\alpha$ asserts that the context Γ entails that $a:\alpha$. The notation $\Gamma, J \vdash J'$ abbreviates $\Gamma \cup \{J\} \vdash J'$.

With judgements, we now can state *typing inferences rules*. Such inference rules have the form

$$\frac{P_1 \quad \cdots \quad P_n}{Q},$$

which asserts that the premises P_1, \dots, P_n entail the conclusion Q . For example,

$$\frac{\Gamma \text{ is a judgement context} \quad a \text{ is a term} \quad \alpha \text{ is a type} \quad \Gamma \vdash a:\alpha}{\Gamma \vdash a:\alpha}$$



is a particularly uninteresting inference rule. Explicitly stating the domain of each variable as premises is cumbersome however, so the following are conventions for variable domains based on their name:

► Γ, Γ_i, \dots each range over judgement contexts,

► a, b, c, \dots each range over terms,

► $\alpha, \beta, \gamma, \dots$ each range over types.

The following typing rules are given for Language-A:

$$\begin{array}{lcl}
 \text{TODO} & \frac{}{\Gamma, a:\alpha \vdash a:\alpha} & \text{ } \\
 \text{Function-Abstraction} & \frac{\Gamma, a:\alpha \vdash b:\beta}{\Gamma \vdash a:\alpha \Rightarrow b:\alpha \rightarrow \beta} & \text{ } \\
 \text{Function-Concretization} & \frac{\Gamma \vdash a:\alpha \rightarrow \beta \quad \Gamma \vdash b:\alpha}{\Gamma \vdash a \ b:\beta} & \text{ } \\
 \text{Construct-Sum-Left} & \frac{\Gamma \vdash a:\alpha}{\Gamma \vdash (\text{left } a):\alpha + \beta} & \text{ } \\
 \text{Construct-Sum-Right} & \frac{\Gamma \vdash b:\beta}{\Gamma \vdash (\text{right } b):\alpha + \beta} & \text{ } \\
 \text{Construct-Product} & \frac{\Gamma \vdash a:\alpha \quad \Gamma \vdash b:\beta}{\Gamma \vdash ((a, b)):\alpha \times \beta} & \text{ } \\
 \text{Destruct-Sum} & \frac{\Gamma \vdash x:(\alpha + \beta) \quad \Gamma, a:\alpha \vdash c_1:\gamma \quad \Gamma, b:\beta \vdash c_2:\gamma}{\Gamma \vdash (\text{match } x \{ \text{left } a \Rightarrow c_1 \mid \text{right } b \Rightarrow c_2 \}): \gamma} & \text{ } \\
 \text{Destruct-Product-First} & \frac{\Gamma \vdash x:(\alpha \times \beta)}{\Gamma \vdash (\text{first } x):\alpha} & \text{ }
 \end{array} \tag{1.2}$$

Reduction Rules for Language-A

Finally, the last step is to introduce *reduction rules*. So far we have outlined syntax and inference rules for building expressions in Language-A, but all these expressions are inert. Reduction rules describe how terms can be transformed, step by step, in a way that models computation. A series of these simple reductions may end in a term for which no reduction rule can apply. Call these terms *values*, and notate “ v is a value” as “value v .” The following reduction rules are given for Language-A:

$$\beta\text{-Reduction} \quad \frac{\Gamma \vdash (a:\alpha \Rightarrow b):(\alpha \rightarrow \beta) \quad \Gamma \vdash v:\alpha \quad \text{value } v}{(a:\alpha \Rightarrow b) \ v \rightarrow [v/a]b} \tag{1.3}$$

The most fundamental of these rules is β -Reduction, which is the way that function applications are resolved to the represented computation’s output. The substitution notation $[v/a]b$ indicates to “replace with v each appearance of a in b .” In this way, for a function $a:\alpha \Rightarrow b:(\alpha \rightarrow \beta)$ and an input $v:\alpha$, β -Reduction *substitutes* the input v for the appearances of the function parameter a in the function body b .

For example, consider the following terms: **TODO**: enlightening example of a series of β -reductions for some simple computation.

1.2.2 Other Structures in Language-A

TODO: decide, based on future developments, which terms and types should be included here. Things like:

- integer
- natural
- boolean
- unit

1.2.3 Properties of Language-A

With the syntax, typing rules, and reduction rules for Language-A, we now have a completed definition of the language. However, some of the design decisions may seem arbitrary even if intuitive. This particular framework is good because it maintains a few nice properties that make reasoning about Language-A intuitive and extendable.

TODO: how much detail do I want to go into this? come back after writing rest of chapter 1.

Type-Preserving Substitution

$$\Gamma, a:\alpha \vdash b:\beta \wedge \Gamma \vdash v:\alpha \wedge \text{value } v \implies \Gamma \vdash [v/a]b:\beta \quad (1.4)$$

Reduction Progress

$$\{\} \vdash a:\alpha \implies \text{value } a \vee \exists a':a \twoheadrightarrow a' \quad (1.5)$$

Type Soundness

$$\Gamma \vdash a:\alpha \wedge a \twoheadrightarrow a' \implies \text{value } a' \vee \exists a'':\alpha \twoheadrightarrow a'' \quad (1.6)$$

Type Preservation

$$\Gamma \vdash a:\alpha \wedge a \twoheadrightarrow a' \implies \Gamma \vdash a':\alpha \quad (1.7)$$

1.3 Computation with Effects

The definition of Language-A in section ?? gives a good general flavor for many similar functional programming languages. In terms of the such language's reductions from term to term, all the information relevant for deciding such reductions is explicit within the term itself and the explicit context built up during reduction. In other words, there is no *implicit activity* that influences what a term's reduction will look like. The path of reductions is exactly the computation a term corresponds to, so this yields that the computation has the same property of not having access to or being affected by any implicit activity.

However nice a formalization of computation this is, it is immediately unrealistic. Actual computers, for which programming languages are abstractions of their activities, host multitudes of implicit processes while running a program. Even if a programming language modelled all such processes and incorporated them into the language so that each activity was made perfectly explicit as a term (a task that is certainly infeasible and likely impossible), the result would be an unuseful and inefficient language. The point of having layers of abstraction, in the form of high-and-higher level programming language, is to avoid this situation in the first place. So there appears to be a dilemma:

The Dilemma of Implicit Activities

- (1) Require fully explicit terms and reductions. This grants reasoning about programs is fully formalized and abstracted from the annoyances of hardware and lower-level-implementations, but restricts such programs from being applicable in almost all useful circumstances.
- (2) Allow implicit activities that affect reductions. This grants many useful program applications and maintains some formal nature to the language's behavior, but reasoning about programs is now inescapably tainted by implicit activities.

A resolution to this apparent dilemma is to either choose one horn or to reject the dilemma. But first, let us consider what kinds of implicit activities are being considered here — in particular, what are computational *effects*.

1.3.1 Effects

The definition of an *effect* in the context of programming language is frequently debated, and there is no majority standard answer¹. For the purposes of this thesis, the following definition is adopted.

Definition 1.3.1. A computational *effect* is a capability in a program that depends on factors outside of that capability's normal scope.

¹**TODO:** source

The definition of *normal scope* will be left to intuitive interpretation, with the intent that what is the normal scope of a capability depends on many theoretical, design, and implementation factors.

Now for example, suppose we have a program P that computes the sum of two integers given as input. If P is designed and implemented exactly to this specification, then P has no effects; none of P 's capabilities depend on factors outside of their usual scope. The only scope that P 's capabilities depend on is that of the two inputs. No other factors influence what the correctly computed sum of the two integers is. Such a program with no effects is called a *pure* program.

Definition 1.3.2. A program is *pure* if it has no effects.

But if P is run, as abstractly as it is defined, not much use comes from it. P does not display its results, write the results to some P -specified memory, or give you an error if there is an overflow. These capabilities depend on factors outside of P 's normal scope, per its specification, and so are examples of effects.

So as another example, let us consider a modified version of P that is effectual. Suppose that program P' takes as input two integers, computes the sum of the integers, throws an error if there is an overflow, writes the result into RAM, and prints the result to the console from which P' was run. These capabilities are examples of effects, since their behavior depends on factors outside of the normal scope of P' (the same scope as P). Such a program with effects is called an *impure* program.

Definition 1.3.3. A program is *impure* if it has effects.

There are a variety of common effects that are available in almost every programming language. These include:

- *input/output (IO)*, e.g. printing to the console, accepting input from user, interfacing with other peripherals.
- *mutable data*, e.g. mutable variables, in-place arrays.
- *exception*, e.g. division by 0, out-of-bounds index of array.
- *nondeterminism*, e.g. **TODO**: what would be good examples for this...
- *partiality*, e.g. **TODO**
- *continuation*, e.g. **TODO**

TODO: go into detail here, or till after I've talked about comparing functional/imperative approaches to effects?

According to definition 1.3.1, there is an easy method for transforming an impure program into a pure program: add the factors depended upon by the program's effects to the program's scope. Using this intuition, we can reformulate the Dilemma of Implicit Activities with the new formal notion of computational effects:

The Dilemma of Scopes for Effects

- (1) Require only pure programs. This grants reasoning about programs to depend only on the normal scope of the program.
- (2) Allow impure as well as pure programs. This grants many useful programs, where the behavior of the programs depends on factors not entirely encapsulated by the program's normal scope.

The goal of resolving this dilemma is to make a choice about how programming languages should be designed. Are computational effects a feature to be avoided as much as possible? Or are they actually so necessary that it would be a mistake to restrict them? Unsurprisingly, no widely-adopted languages have chosen horn (1). But even in more generality, almost all languages have chosen horn (2), but with many varied approaches to incorporating effects. Overallly languages have clustered into two groups.

- *Functional* programming language treat computation as the evaluation of mathematical functions. Such languages put varying degrees of emphasis on restricting effects, but in general much more emphasis than imperative languages.
- *Imperative* programming languages treat computation as a sequence of commands. E.g. the C family, Bash, Java, Python. Such languages put very little (if any) emphasis on restricting effects. E.g. Scheme, Lisp, SML, Clojure, Scala, Haskell, Agda, Gallina.

The perspective on effects is not the only way in which these groups differ, but nevertheless it is an important one. Among functional programming languages, we shall consider two in particular: SML and Haskell. SML takes a very relaxed approach to effects, and Haskell takes a more restrictive approach.



1.3.2 A Simple Approach to Effects

The SML (Standard Meta Language) language is most commonly used in academic settings for teaching about programming languages. However here, I shall avoid introducing an entire new syntax and language semantics. Instead, the SML effect impkementation strategy will be demonstrated using Language-B, a variant of Language-A with the following new features:

- *Sequence*. A collection of terms is arranged in the sequence of desired resolution. Reducing the sequence term will reduce the terms in order.
- *Mutable*. A new type that indicates mutable data. It must be explicitly read from and wrote to, rather than handled like the normal (immutable) value.
- *Exception*. A collection of new term constructs, allowing for programs to *throw* and *catch* exceptions that interrupt usual reduction.

Syntax for Language-B

metavariable	constructors	name
$\langle term \rangle$	$\langle term \rangle ; \langle term \rangle$ $@\langle term \rangle$ $!\langle term \rangle$ $\langle term-name \rangle := \langle term \rangle$ $throw \langle term \rangle$ $catch \langle term \rangle \{ \langle name \rangle \Rightarrow \langle term \rangle \}$	sequence mutable read mutable write mutable throw exception catch exception
$\langle type \rangle$	$mutable \langle type \rangle$	mutable

(1.8)

Typing Rules for Language-B

Construct-Sequence	$\frac{\Gamma \vdash a:\alpha \quad \Gamma \vdash b:\beta}{\Gamma \vdash (a ; b):\beta}$
Construct-Mutable	$\frac{\Gamma \vdash a:\alpha}{\Gamma \vdash (@a):mutable \alpha}$
Construct-Mutable-Read	$\frac{\Gamma \vdash m:mutable \alpha}{\Gamma \vdash (!m):\alpha}$
Construct-Mutable-Write	$\frac{\Gamma \vdash m:mutable \alpha \quad \Gamma \vdash a:\alpha}{\Gamma \vdash (m := a):unit}$
Construct-Throw	$\frac{}{(throw e):\epsilon}$
Construct-Catch	$\frac{\Gamma \vdash a:\alpha \quad \Gamma, e:\epsilon \vdash a':\alpha}{(catch a \{ e \Rightarrow a' \}):\alpha}$

(1.9)

Reduction Rules for Language-B

TODO: formally explain how environments work (E). Should that go in this section, or the definition of Language-A? If I put it in Language-A, that would make defining let expressions much easier.

TODO: formally explain how mutable environments work (M). Models the storing of values for mutables.

$$\begin{array}{ll}
\text{Sequence-Reduce} & \frac{a \rightarrow a'}{(a ; b) \rightarrow (a' ; b)} \\
\\
\text{Truncate} & \frac{\text{value } a}{(a ; b) \rightarrow b} \\
\\
\text{Initialize} & \frac{\text{value } v}{@v \rightarrow \bullet \quad M[@v] = v} \\
\\
\text{Mutable-Reduce} & \frac{a \rightarrow a'}{@a \rightarrow @a'} \\
\\
\text{Read} & \frac{m \in M}{!m \rightarrow M[m]} \tag{1.10} \\
\\
\text{Read-Reduce} & \frac{m \rightarrow m'}{!m \rightarrow !m'} \\
\\
\text{Write} & \frac{m \in M \quad \text{value } v}{m := v} \\
\\
\text{Write-Source-Reduce} & \frac{m \rightarrow m'}{m := a \rightarrow m' := a} \\
\\
\text{Write-Target-Reduce} & \frac{a \rightarrow a'}{m := a \rightarrow m := a'}
\end{array}$$

TODO: reduction rules for throw/catch.