

Aspects of efficiency in functional programming languages

by

Samuel Valdemar Grange

supervised by

Prof. Kim Skak Larsen



UNIVERSITY OF SOUTHERN DENMARK
DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE
Master's thesis in Computer Science

Contents

I	Compilers and languages	2
1	Programming languages	3
1.1	Untyped lambda calculus	4
1.2	Translation to lambda calculus	5
1.2.1	Scoping	5
1.2.2	Recursion	6
1.3	High level abstractions	8
1.3.1	Algebraic data types	8
2	Typing and validation	11
2.1	Types and validation	11
2.1.1	The language of types	11
2.2	Hindley-Milner	14
2.2.1	Damas-Milner Algorithm W	17
2.2.2	Recursion	19
2.3	The cost of expressiveness	19

Part I

Compilers and languages

Chapter 1

Programming languages

Computers are devices which read a well-defined, finite sequence of simple instructions and emit a result. In theoretical analysis of computers, models have been developed to understand and prove properties. A finite sequence of instructions fed to a computer is called an *algorithm*, which is the language of high level computation [3]. In modern encodings of algorithms or programs, “high level” languages are used instead of the computational models. Such languages are then translated into instructions that often are much closer to a computational model. The process of translating programs into computer instructions is called *compiling*, or *transpiling* if the program is first translated into another “high level” language.

For the purpose of this dissertation, a simple programming language has been implemented to illustrate the concepts in detail. The language transpiles to *untyped lambda calculus*. For the remainder, the language will be referred to as *L*.

1.1 Untyped lambda calculus

The *untyped lambda calculus* is a model of computation, developed by Alonzo Church[1]. The untyped lambda calculus is a simple tangible language of just three terms.

$$x \tag{1.1}$$

$$\lambda x.E \tag{1.2}$$

$$YE \tag{1.3}$$

The first term is that of a *variable* (Equation 1.3). A variable is a reference to another lambda abstraction. Equation 1.2 shows a lambda *abstraction*, which contains a *bound* variable x and another lambda term E . A bound variable is one explicitly parameterized in the abstraction. In traditional programming languages bound variables are parameters to a function, whereas *free* variables are variables that exist outside of the function scope. It is easy to see that variables are either free or bound. Free variables are determined by Equation 1.4, Equation 1.5 and Equation 1.6.

$$Free(x) = \{x\} \tag{1.4}$$

$$Free(\lambda x.E) = Free(E) \setminus \{x\} \tag{1.5}$$

$$Free(YE) = Free(Y) \cup Free(E) \tag{1.6}$$

Finally, in Equation 1.3, *application*. An application of two terms can be interpreted as substituting the variable in the left abstraction Y with the right term E . It is also common to introduce the *let binding* to lambda calculus, which will be done when introducing typing in ??.

Example 1.1.1. Let Y be $\lambda x.T$ and E be z , then YE is $(\lambda x.T)z$. Furthermore, substituting x for E , such that Y becomes $T[x := E]$. Since $E = z$, substitute E for z , such that $T[x := z]$, read as “Every instance of x in T , should be substituted by z ”.

Remark 1.1.1. Substituting lambda terms is a popular method of evaluating lambda calculus programs. Languages like Miranda, Clean and general purpose evaluation programs like the G-machine ??, implement *graph rewriting*, which will be introduced in ??

The untyped lambda calculus is in fact turing complete; any algorithm that can be evaluated by a computer can be encoded in the untyped lambda calculus. The turing completeness of the untyped lambda calculus can be realized by modelling numerics, boolean logic and recursion with the *Y-combinator*. Church encoding, is the encoding of numerics, arithmetic expressions and boolean logic [2]. For the remainder of the dissertation, ordinary arithmetic expressions are written in traditional mathematics. The expressiveness and simplicity of lambda calculus, makes it an excellent language to transpile to, which in fact, is a common technique.

1.2 Translation to lambda calculus

High level languages associated with lambda calculus are often also very close to it. The L language is very close to the untyped lambda calculus. See two equivalent programs, Equation 1.7 and Listing 1.1, that both add an a and a b .

$$(\lambda add.E)(\lambda a.(\lambda b.a + b)) \tag{1.7}$$

Listing 1.1: Add function

```
1 fun add a b = a + b;
```

Notice that in Equation 1.7 the term E is left undefined, E is “the rest of the program in this scope”. If the program was to apply 1 and 2 to add, directly below in the high level representation (Listing 1.2) the lambda calculus equivalent would look like Equation 1.8.

$$(\lambda add.add\ 1\ 2)(\lambda a.(\lambda b.a + b)) \tag{1.8}$$

Listing 1.2: Add function applied

```
1 fun add a b = a + b;  
2 add 1 2;
```

1.2.1 Scoping

Notice that Equation 1.7, must bind the function name “outside the rest of the program” or more formally in an outer scope. In a traditional program such as Listing 1.3, functions must be explicitly named to translate as in the above example.

Listing 1.3: A traditional program

```
1 fun add a b = a + b;  
2 fun sub a b = a - b;  
3 sub (add 10 20) 5;
```

Listing 1.4: An order dependant program

```
1 fun sub a b = add a (0 - b);  
2 fun add a b = a + b;  
3 sub (add 10 20) 5;
```

Notice that there are several problems, such as, the order of which functions are defined may alter whether the program is correct or not. For instance, the program defined in Listing 1.4 would not translate correct, it would translate to Equation 1.9. The definition of *sub*, or rather, the applied lambda abstraction, is missing a reference to the *add* function.

$$(\lambda sub.(\lambda add.(sub (add 10 20) 5)) (\lambda a.(\lambda b.a + b))) (\lambda a.(\lambda b.add a(0 - b))) \quad (1.9)$$

lambda lifting is a technique where free variables (section 1.1), are explicitly parameterized [9]. This is exactly the problem in Equation 1.9, which has the lambda lifted solution Equation 1.10.

$$(\lambda sub.(\lambda add.(sub add (add 10 20) 5)) (\lambda a.(\lambda b.a + b))) (\lambda add.(\lambda a.(\lambda b.add a(0 - b)))) \quad (1.10)$$

As it will turn out, this will also enables complicated behaviour, such as *mutual recursion*.

Moreover, lambda lifting also conforms to “traditional” scoping rules. *Variable shadowing* occurs when there exists $1 < \text{reachable variables of the same name}$, but the “nearest”, in regard to scope distance is chosen. Effectively, other variables than the one chosen, are *shadowed*. Variable shadowing is an implied side-effect of using using lambda calculus. Convince yourself that the function *f* in Listing 1.5, yields 12.

Listing 1.5: Scoping rules in programming languages

```

1 let x = 22;
2 let a = 10;
3 fun f =
4   let x = 2;
5   a + x;
```

1.2.2 Recursion

Reductions in mathematics and computer science are one of the principal methods used for developing beautiful equations and algorithms.

Listing 1.6: Infinite program

```

1 fun f n =
2   if (n == 0) n
3   else if (n == 1) n + (n - 1)
4   else if (n == 2) n + ((n - 1) + (n - 2))
5   ...
```

Listing 1.6 defines a function f , that in fact is infinite. Looking at the untyped lambda calculus, there are not any of the three term types that define infinite functions or abstractions, at first glance. Instead of writing an infinite function, the question is rather, how can a reduction be performed on this function, such that it can evaluate *any* case of n ?

Listing 1.7: Recursive program

```

1 fun f n =
2   if (n == 0) n
3   else n + (f (n - 1))

```

Listing 1.7 defines a recursive variant of f , it is a product of the reduction in Equation 1.11.

$$n + (n - 1) \cdot \dots + 0 = \sum_{k=0}^n k \quad (1.11)$$

But since the untyped lambda calculus is turing complete, or rather, if one were to show it were, it must also realize algorithms that are recursive or include loops, the two of which are equivalent in expressiveness.

$$(\lambda f.E)(\lambda n.\text{if}(n == 0)(n)(n + (f(n - 1)))) \quad (1.12)$$

The naive implementation of a recursive variant, will yield an unsolvable problem, in fact, an infinite problem. In Equation 1.12, when f is applied recursively, it must be referenced. How can f be referenced, if it is “being constructed”? Substituting f with its implementation in Equation 1.13, will yield the same problem again, but at one level deeper.

$$(\lambda f.E)(\lambda n.\text{if}(n == 0)(n)(n + ((\lambda n.\text{if}(n == 0)(n)(n + (f(n - 1))))(n - 1)))) \quad (1.13)$$

One could say, that the problem is now recursive. Recall that lambda lifting (subsection 1.2.1), is the technique of explicitly parameterizing outside references. Convince yourself that f lives in the scope above its own body, such that, when referencing f from within f , f should be parameterized as in Listing 1.8, translating to Equation 1.14.

Listing 1.8: Explicitly passing recursive function

```

1 fun f f n =
2   if (n == 0) n
3   else n + (f f (n - 1))

```

$$(\lambda f.E)(\lambda f.(\lambda n.\text{if}(n == 0)(n)(n + (f f (n - 1))))) \quad (1.14)$$

The initial invocation of f , must involve f , such that it becomes $f f n$. The *Y-combinator*, an implementation of a fixed-point combinator, in Equation 1.15 is the key to realize that the untyped lambda calculus can implement recursion. Languages with functions and support binding functions to parameters, can implement recursion with the Y-combinator.

$$\lambda f.(\lambda x.f(xx))(\lambda x.f(xx)) \quad (1.15)$$

Implementing mutual recursion is an interesting case of lambda lifting and recursion in untyped lambda calculus.

Listing 1.9: Mutual recursion

```
1 fun g = f;
2 fun f = g;
```

Notice in Listing 1.9 that g requires f to be lifted and f requires g to be lifted. If a translation “pessimistically” lifts all definitions from the above scope, then all required references exist in lexical scope.

Languages have different methods of introducing recursion, some of which have very different implications, especially when considering types. For instance, OCaml has the `let rec` binding, to introduce recursive definitions. The `rec` keyword indicates to the compiler that the binding should be able to “see itself” (??).

1.3 High level abstractions

The lambda calculus is a powerful language that can express any algorithm. Expressiveness does not necessarily imply ergonomics or elegance, in fact encoding moderately complicated algorithms in lambda calculus becomes quite messy. Many high level techniques exist to model abstractions in tangible concepts.

1.3.1 Algebraic data types

Algebraic data types are essentially a combination of disjoint unions, tuples and records. Algebraic data types are closely related to types thus require some type theory to fully grasp. Types are explored more in depth in ??.

Listing 1.10: List algebraic data type

```
1 type List a =
2   | Nil
3   | Cons a (List a)
4 ;
```

Listing 1.10 is an implementation of a linked list. The list value can either take the type of `Nil` indicating an empty list, or it can take the type of `Cons` indicating a pair of type a and another list. The list implementation has two constructors and one type parameter. The type parameter a of the list algebraic data type defines a *polymorphic type*; a can agree on any type, it is universally quantified $\forall a. \text{Cons } a \text{ (List } a)$. The two constructors `Nil` and `Cons` both create a value of type `List a` once instantiated.

Listing 1.11: List instance and match

```

1 let l = Cons 1 (Cons 2 (Cons 3 Nil));
2
3 match l
4   | Nil -> 0;
5   | Cons n _ -> n;
6 ;

```

Once a value is embedded into an algebraic data type such as a list it must be extractable to be of any use. Values of algebraic data types are extracted and analysed with *pattern matching*. Pattern match comes in many forms, notably it allows one to define a computation based on the type an algebraic data type instance realizes (Listing 1.11).

Scott encoding

Pattern matching strays far from the simple untyped lambda calculus, but can in fact be encoded into it. The *scott encoding* (Equation 1.16) is a technique that describes a general purpose framework to encode algebraic data types into lambda calculus [11]. Considering an algebraic data type instance as a function which accepts a set of “handlers” allows the encoding into lambda calculus. The scott encoding specifies that constructors should now be functions that are each parameterized by the constructor parameters $x_1 \dots x_{A_i}$ where A_i is the arity of the constructor i . Additionally each of the constructor functions return a n arity function, where n is the cardinality of the set of constructors. Of the n functions, the constructor parameters $x_1 \dots x_{A_i}$ are applied to the i ’th “handler” c_i . These encoding rules ensure that the “handler” functions are provided uniformly to all instances of the algebraic data type.

$$\lambda x_1 \dots x_{A_i}. \lambda c_1 \dots c_n. c_i x_1 \dots x_{A_i} \quad (1.16)$$

Example 1.3.1. The `List` algebraic data type in Listing 1.10 has two constructors, `Nil` with the constructor type Equation 1.17 and `Cons` with the constructor type Equation 1.18. Equation 1.17 is in fact also the type of `List` once instantiated, effectively treating partially applied functions as

data.

$$b \rightarrow (a \rightarrow \text{List } a \rightarrow b) \rightarrow b \quad (1.17)$$

$$(a \rightarrow \text{List } a \rightarrow b) \rightarrow b \rightarrow (a \rightarrow \text{List } a \rightarrow b) \rightarrow b \quad (1.18)$$

Listing 1.12: List algebraic data type implementation

```

1 fun cons x xs =
2   fun c _ onCons = onCons x xs;
3   c;
4
5 fun nil =
6   fun c onNil _ = onNil;
7   c;

```

Encoding the constructors in L yields the functions defined in Listing 1.12. Pattern matching is but a matter of applying the appropriate handlers. In Listing 1.13.

Listing 1.13: Example of scott encoded list algebraic data type

```

1 let l = cons 1 (cons 2 (cons 3 nil));
2
3 fun head x _ =
4   x;
5
6 l head 0;

```

Efficiency can be a bit tricky in lambda calculus as it is at the mercy of implementation. A common method of considering efficiency is counting β -reduction since they evaluate to function invocations. The β -reduction is a substitution which substitutes an application where the left side is an abstraction in which the bound variable is substituted with the right side term (Equation 1.19).

$$\beta_{red}((\lambda x.T)E) = T[x := E] \quad (1.19)$$

It should be clear that invoking a n arity function will take n applications. In the case of a scott encoded algebraic data types the largest term in regard to complexity is either the size of the set of “handler” functions or the “handler” function with most parameters. The time to evaluate pattern match is thus $O(\max_i(c_i + A_i))$.

Chapter 2

Typing and validation

Automatic validation is one of many reasons to use computers for solving various tasks including writing new computer programs. Spellchecking is a common and trivial instance of an input validation algorithm.

2.1 Types and validation

The spell checking equivalent for computer programs could be type checking; a subproblem of validating a programmer’s intuition of a program’s intent. Types also have other properties than simply validating they can in fact be related to theorems to which an implementation is the proof [8].

Listing 2.1: Head implementation

```
1 fun head l: List a → a =  
2   match l  
3     | Cons x _ → x;  
4     | Nil → ?;  
5   ;
```

For instance, consider the implementation of the function with type `List a → a` in Listing 2.1. A total implementation of the function cannot exist.

The type system for the *L* language will be the Hindley-Milner type system [7, 10].

2.1.1 The language of types

Before delving into types, the lambda calculus defined in section 1.1 must be augmented with the *let expression* (Equation 2.1).

$$\text{let } x = Y \text{ in } E \tag{2.1}$$

It should be noted that the let binding can be expressed by abstraction and application (Equation 2.2).

$$(\lambda x.E)(Y) \quad (2.2)$$

The let expression has a nice property that will become apparent later when typing rules are introduced.

Types are an artificial layer atop of a program just as spell checking is an artificial layer atop text. There are two variants of types in the Hindley-Milner type system, the *monotype* and the *polytype*. A monotype is either a type variable, an abstraction of two monotypes or an application of a type constructor (Equation 2.3).

$$mono \ \tau = a \mid \tau \rightarrow \tau \mid C\tau_1 \dots \tau_n \quad (2.3)$$

Atoms are terminal terms in a formula and are expressed either by type variable a or C with no type parameters. The application term of the monotype is dependent on the primitive types of the programming language. The types $\tau_1 \dots \tau_n$ are monotype parameters required to construct some type C . In L the set of type constructors are $\{\mathbf{Int}, \mathbf{Bool}\} \cup \mathbf{ADT}$. \mathbf{Int} and \mathbf{Bool} are type constructors of arity 0 thus only have one instantiation and are atomic. The set of constructors \mathbf{ADT} encapsulates the set of program defined algebraic data structures (??).

Example 2.1.1. Let $\mathbf{ADT} = \{\mathbf{List}\}$ where \mathbf{List} is defined as in Listing 1.10. The *type constructor* (not to be confused for constructors like \mathbf{Cons} or \mathbf{Nil}) for \mathbf{List} has the signature $\mathbf{a} \rightarrow \mathbf{List} \ \mathbf{a}$ stating that if supplied with some type \mathbf{a} it constructs a type of $\mathbf{List} \ \mathbf{a}$ (effectively containing the provided type). The type \mathbf{List} is a type constructor with one type parameter \mathbf{a} .

\perp denotes falsity, in type systems a value of this type can never exist since that in itself would disprove the program. It is common in programming languages with strong type systems to let thrown exceptions be of type \perp since it adheres to every type and indicates that the program is no longer running, since no instance of \perp can exist. \top denotes truth, in type systems every type is a supertype of \top . \top is in practice only used to model side effects, since not all side effects return useful values. In programming languages with side effects \perp and \top are considerably more useful than in pure programming languages.

A polytype is a polymorphic type (Equation 2.4).

$$poly \ \sigma = \tau \mid \forall a. \sigma \quad (2.4)$$

Polymorphic types either take the shape of a type variable or introduce a type which all types a adhere to. This does not necessarily include *all* types since the **Gen** rule of Figure 2.2 constrains the domain that a ranges over to contain only bound type variables.

Remark 2.1.1. All types adhere to a polymorphic type but polymorphic types do not adhere to any type other than polymorphic types. The concept of adherence in types is commonly called *subtyping*. Every subtype is a *at least* an implementation of it's supertype. Since this concept can be difficult to grasp from just text, observe Figure 2.1. Note that σ is controversial to

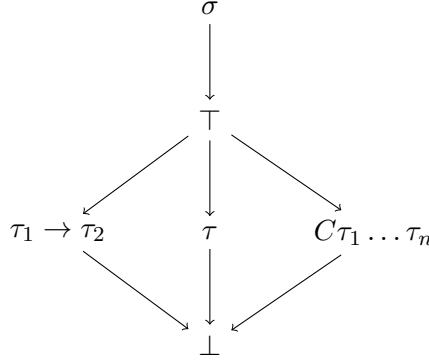


Figure 2.1: The type hierarchy of Hindley-Milner.

introduce to the type hierarchy and has only been so to illustrate the point of subtyping. σ is but a mechanism to prove type systems, σ is never a specific type.

A principal component of typing in Hindley-Milner is the *environment*. The environment Γ is a set of pairs of variable and type (Equation 2.5). $\Gamma \vdash x : \sigma$ implies a *typing judgment*, meaning that given Γ , the variable x can adhere to the type σ .

Remark 2.1.2. Judging a type does not necessarily mean that the judged type is the only type that x may take, it states that it is one *possible* type that x may take. The property of taking multiple possible types is what allows polymorphism. This is made more apparent in Example 2.2.2 where `id` may take the type of either $\forall a. a \rightarrow a$, $\text{Int} \rightarrow \text{Int}$ or $\forall a. (a \rightarrow a) \rightarrow (a \rightarrow a)$.

$$\Gamma = \epsilon \mid \Gamma, x : \sigma \quad (2.5)$$

Like in the untyped lambda calculus, types also have notions of free and bound type variables. Type variables are bound when they have been

introduced by a quantification or exist in the environment.

$$\text{free}(a) = \{a\} \quad (2.6)$$

$$\text{free}(C\tau_1 \dots \tau_n) = \bigcup_{i=1}^n \text{free}(\tau_i) \quad (2.7)$$

$$\text{free}(\tau_1 \rightarrow \tau_2) = \text{free}(\tau_1) \cup \text{free}(\tau_2) \quad (2.8)$$

$$\text{free}(\Gamma) = \bigcup_{x:\sigma \in \Gamma} \text{free}(\sigma) \quad (2.9)$$

$$\text{free}(\forall a.\sigma) = \text{free}(\sigma) - \{a\} \quad (2.10)$$

Free and bound type variables are a vital distinction for type analysis since **Gen** in Figure 2.2 may only generalize over bound type variables. The variables which may appear in a quantification have an important role in Equation 2.11, since only free variables may be substituted.

Example 2.1.2. Consider the type for the function **fst** in Listing 2.3.

Listing 2.2: First function

```
1 fun fst a b: ∀A.∀B.A → B → A = a;
```

Listing 2.3: First function in lambda calculus

```
1 let fst = λa.(let f = λb.a in f) in fst
```

The type for **fst** is $\forall A \forall B. A \rightarrow B \rightarrow A$.

Note that a naive typing would look something like $\forall A. A \rightarrow (\forall B. B \rightarrow A)$ but rank-2 polymorphism is not typable in Hindley-Milner. An important realization is the context from where the type analysis is made. If type analysis is made from within **f** the type of **f** becomes $\forall B. B \rightarrow A$ and the type variable A is free.

2.2 Hindley-Milner

With the now introduced primitives, the Hindley-Milner type system is but a set of inference rules composed by said primitives. There are six rules in the Hindley-Milner rules outlined in Figure 2.2.

- **Var** states that if some variable x with type σ exists in the environment, the type can be judged. In practice, when $x : \sigma$ is encountered in the expression tree it is added to the environment.
- **App** decides that if $e_1 : \tau_1 \rightarrow \tau_2$ and $e_2 : \tau_1$ has been judged to exist then $e_1 e_2$ implies the removal of τ_1 from $\tau_1 \rightarrow \tau_2$ such that $e_1 e_2 : \tau_2$.

$\text{Var} \frac{x : \sigma \in \Gamma}{\Gamma \vdash x : \sigma}$	$\text{App} \frac{\Gamma \vdash e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash e_2 : \tau_1}{\Gamma \vdash e_1 e_2 : \tau_2}$
$\text{Abs} \frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x. e : \tau_1 \rightarrow \tau_2}$	$\text{Let} \frac{\Gamma \vdash e_1 : \sigma \quad \Gamma, x : \sigma \vdash e_2 : \tau}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \tau}$
$\text{Ins} \frac{\Gamma \vdash e : \sigma_1 \quad \sigma_1 \sqsubseteq \sigma_2}{\Gamma \vdash e : \sigma_2}$	$\text{Gen} \frac{\Gamma \vdash e : \sigma \quad a \notin \text{free}(\Gamma)}{\Gamma \vdash e : \forall a. \sigma}$

Figure 2.2: Hindley-Milner type rules

- **Abs** is the typing rule of lambda abstractions. If $x : \tau_1$ exists in the environment from some type analysis of e and the abstraction's body e has been judged to be of type τ_2 then the abstraction of x must take the type of x to create the type of the body e .
- **Let** states that if e_1 has been judged to have type σ then the let expression's identifier $x : \sigma$ must exist in the environment when deriving the type of e_2 . Observe that **Let** introduces a polymorphic type to the environment while **Abs** introduces a monomorphic one. Note that by Remark 2.1.2 x may be polymorphic in e_2 .
- **Inst** specializes some polymorphic type (in regard to the type system implementation) to a more specific polymorphic type. \sqsubseteq is the partial order of types where the binary relation between two types compares the descriptiveness of types.

Example 2.2.1. In L the smallest element is the top of the type hierarchy (Figure 2.1), the polymorphic type.

- **Gen** generalizes over all bound variables a .

Let polymorphism is exemplified in Example 2.2.2.

Example 2.2.2. Throughout this example (x, z) is the pair of the variables x and z .

The identity function is a common example to illustrate type systems (Listing 2.4).

Listing 2.4: Identity function in L

```

1 fun id x = x;
2 id 4;
```


Listing 2.5: Identity function in lambda calculus with let

```
1 let id = (λx.x) in
2 id 4
```

Stating that `id` has the type $\forall a.a \rightarrow a$ and `4` has the type `Int` is Listing 2.5 program correct? By applying the Hindley-Milner rules one can prove or disprove this statement. A correct proof of Listing 2.4 must be Figure 2.3.

Listing 2.6: Identity function in lambda calculus by abstraction

```
1 (λid.id 4)(λx.x)
```

Listing 2.5 and Listing 2.6 are two equivalent programs with slightly different proofs which raises the question of why the `let` expression is even needed. If Listing 2.5 and Listing 2.6 were to be slightly changed such that two new programs Listing 2.7 and Listing 2.8 were to be proved, Listing 2.8 would not be provable while Listing 2.7 would.

Listing 2.7: Identity function with two applications

```
1 let id = (λx.x) in
2 (id 4, id id)
```

Listing 2.8: Identity function with two applications as abstraction

```
1 (λid.(id 4, id id)(λx.x)
```

In Listing 2.8 `id` cannot adhere to polymorphism by **Abs** in Figure 2.2 whilst **Let** can.

$$\begin{array}{c}
 \text{Var} \frac{id : \forall a.a \rightarrow a \in \Gamma}{\Gamma \vdash id : \forall a.a \rightarrow a} \quad \frac{\forall a.a \rightarrow a \sqsubseteq \text{Int} \rightarrow \text{Int}}{\Gamma \vdash id : \text{Int} \rightarrow \text{Int}} \quad \frac{4 : \text{Int} \in \Gamma}{\Gamma \vdash 4 : \text{Int}} \text{Var} \\
 \text{Inst} \frac{\Gamma \vdash id : \forall a.a \rightarrow a \quad \Gamma \vdash id : \text{Int} \rightarrow \text{Int} \quad \Gamma \vdash 4 : \text{Int}}{\Gamma, id : \forall a.a \rightarrow a \vdash id\ 4 : \text{Int}} \text{App} \\
 \\
 \text{(a)} \\
 \\
 \text{Abs} \frac{\text{Var} \frac{x : a \in \Gamma}{\Gamma \vdash x : a}}{\Gamma \vdash (\lambda x.x) : a \rightarrow a} \quad a \notin free(\Gamma) \\
 \text{Gen} \frac{\Gamma \vdash (\lambda x.x) : a \rightarrow a}{\Gamma \vdash (\lambda x.x) : \forall a.a \rightarrow a} \\
 \text{(b)} \\
 \\
 \text{Let} \frac{\text{2.3b} \quad \text{2.3a}}{\Gamma \vdash \text{let id} = (\lambda x.x) \text{ in id 4 : Int}} \\
 \text{(c)}
 \end{array}$$

Figure 2.3: Identity function instantiation proof

$S\Gamma = \{(x, S\sigma) \mid \forall(x, \sigma) \in \Gamma\}$	(Environment)
$S\sigma = \begin{cases} S\tau & \text{if } \sigma \equiv \tau \\ \{a' \mapsto \tau_1 \mid (a', \tau_1) \in S \mid (a, *) \notin S\}\sigma' & \text{if } \sigma \equiv \forall a. \sigma' \end{cases}$	(Poly)
$S(\tau_1 \rightarrow \tau_2) = S\tau_1 \rightarrow S\tau_2$	(Arrow)
$Sa = \begin{cases} \tau & \text{if } (a, \tau) \in S \\ a & \end{cases}$	(Typevariable)
$SC\tau_1 \dots \tau_n = CS\tau_1 \dots S\tau_n$	(Constructor)

Figure 2.4: Substitution semantics

2.2.1 Damas-Milner Algorithm W

Typing rules are by themselves not that useful since they need all type information declared ahead of checking, inference attempts to guess types such that the rules are satisfied. Type inference is the technique of automatically deriving types, of which there exist many algorithms. One of the most common inference algorithms that produce typings which the Hindley-Milner rules accept is the Damas-Milner Algorithm W inference algorithm [4, 5].

The Damas-Milner Algorithm W rules (Figure 2.6) introduce some new concepts such as *fresh variables*, *most general unifier*, and the *substitution set*. Fresh variables are introduced by picking a variable that has not been picked before from the infinite set $\tau_1, \tau_2 \dots$. Fresh variables are introduced when unknown types are discovered and later unified. The substitution set is a mapping from type variables to types (Equation 2.11).

$$S = \{a_1 \mapsto \tau_1, a_2 \mapsto \tau_2 \dots, a_n \mapsto \tau_n\} \quad (2.11)$$

A substitution written ST where T is an arbitrary component of Hindley-Milner like an environment in which all type variables are substituted (Figure 2.4). Substitution sets can also be combined $S_1 \cdot S_2$ with well defined semantics. The combination of substitution sets is a key component for the correctness of the Damas-Milner inference algorithm.

$$S_1 \cdot S_2 = \{(a \mapsto S_1\tau) \mid (a \mapsto \tau) \in (S_2 \cup S_1)\} \quad (2.12)$$

The \cup operator is distinct in the type variable and is left-biased (Equation 2.13).

$$S_1 \cup S_2 = S_1 \cup \{(a \mapsto \tau) \mid (a \mapsto \tau) \in S_2 \mid (a, *) \notin S_1\} \quad (2.13)$$

Remark 2.2.1. By the substitution set combination operator transitive and circular substitutions cannot occur since type variables in S_1 will inherit all the mappings from S_2 and substitute ones that exist in both. The

$$\begin{array}{c}
\text{Arrow} \frac{S, \{(\tau_1 \rightarrow \tau_2, \gamma_1 \rightarrow \gamma_2)\} \cup T}{S, T \cup \{(\tau_1, \gamma_1), (\tau_2, \gamma_2)\}} \\
\\
\text{Intro} \frac{\tau_1, \tau_2}{\emptyset, \{(\tau_1, \tau_2)\}} \quad \text{Var empty} \frac{S, \{(a, \tau_1)\} \cup T}{S, T} \quad a \equiv \tau_1 \\
\\
\text{Var sub} \frac{S, \{(a, \tau_1)\} \cup T \quad a \notin \text{free}(\tau_1)}{S \cup \{a \mapsto S\tau_1\}, \{a \mapsto S\tau_1\} \cup T} \\
\\
\text{Atom} \frac{S, C_1\tau_1 \dots \tau_n, C_2\gamma_1 \dots \gamma_n \cup T \quad C_1 \equiv C_2}{S, \{(\tau_1, \gamma_1) \dots, (\tau_n, \gamma_n)\} \cup T}
\end{array}$$

Figure 2.5: Rules for most general unification

properties ensured by the combination semantics also induce the property of idempotence. This property is enforced by the Damas-Milner Algorithm W inference rules.

Unification is performed differently based on the context. Unification is performed on monotypes, each of which can take one of three forms (Equation 2.3). Note that the Var rules for most general unifier outlined in Figure 2.5 are commutative.

An important implementation detail which should be noted is that of the polymorphic type. Polymorphic types can be regarded as being a pair of bound types and monotype. This representation is convenient for the **Gen** rule.

Remark 2.2.2. The Damas-Milner algorithm W is the most popular inference algorithm for Hindley-Milner. Though it remains the most popular, it has some interesting competitors. One of which is that of the constraint solver approach which is also used in OCaml [6]. The constraint solver approach is a two phase type inference algorithm. In the first phase the algorithm inspects the expression tree and generates a set of constraints as it goes. After the set of constraints C has been generated it then traverses the constraints and generates type variable substitutions. It is argued that error reporting is significantly easier in such an approach.

In the language L functions **fun** translate to let expressions while **let** translates to abstraction and application.

$$\begin{array}{c}
\text{Var} \frac{x : \sigma \in \Gamma \quad \tau = \text{inst}(\sigma)}{\Gamma \vdash x : \tau, \emptyset} \quad \text{Abs} \frac{\tau_1 = \text{fresh} \quad \Gamma, x : \tau_1 \vdash e : \tau_2, S}{\Gamma \vdash \lambda x. e : S\tau_1 \rightarrow \tau_2, S} \\
\\
\text{App} \\
\frac{\Gamma \vdash e_1 : \tau_1, S_1 \quad \tau_3 = \text{fresh} \quad S_1 \Gamma \vdash e_2 : \tau_2, S_2 \quad S_2 = \text{mgu}(S_2\tau_1, \tau_2 \rightarrow \tau_3)}{\Gamma \vdash e_1 e_2 : S_2\tau_3, S_3 \cdot S_2 \cdot S_1} \\
\\
\text{Let} \frac{\Gamma \vdash e_1 : \tau_1, S_1 \quad S_1 \Gamma, x : S_1\Gamma(\tau_1) \vdash e_2 : \tau_2, S_2}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \tau_2, S_1 \cdot S_2}
\end{array}$$

Figure 2.6: Algorithm W

2.2.2 Recursion

2.3 The cost of expressiveness

Modern languages with strong type systems tend to be notoriously slow to type. In fact many languages with strong type systems provide type systems expressive enough to be turing complete. Though languages with strong type systems can be slow it is most commonly observed from pathological inputs.

In the construction of the compiler for L one target was the C++ language. An instance of a pathological input for the C++ type checker is most definitely the untyped lambda calculus. The lambda terms in C++ must adhere to polymorphism in many cases which has some unknown but large blowup in compilation time. In fact polymorphism is commonly the root of blowup in typing.

ML which implemented Hindley-Milner was believed to have linear complexity before shown to be exponential. As it will turn out, Hindley-Milner also suffers an explosive worst case induced by a pathological input fueled by polymorphism.

Lemma 3. *Hindley-Milner typing is worst case $O(c^n)$.*

Proof. The basis of the blowup stems from the introduced fresh type variables in the polymorphic let inference rule. If the amount of type variables can be shown to be exponential the running time must be at least the same by operations such as subsection set combination and unification.

Listing 3.1: Nested pair

```

1 fun dup a f = f a a;
2 fun deep x = dup (dup (dup (... x)));

```

Listing 3.1 builds a large function signature for **deep**. The innermost **pair** invocation will have its signature unified to $x \rightarrow (x \rightarrow x \rightarrow \tau) \rightarrow \tau$. The second innermost **pair** invocation has the signature $((x \rightarrow x \rightarrow \tau) \rightarrow \tau) \rightarrow (((x \rightarrow x \rightarrow \tau) \rightarrow \tau) \rightarrow \gamma) \rightarrow \gamma$. Naively one might judge Listing 3.1 to run in $O(2^n)$ but an important observation for why Listing 3.1 does not induce exponential blowup is the uniqueness of the type variables. If an efficient representation of **dup** was implemented such that the left and right side were shared such that $\mu \mapsto ((x \rightarrow x \rightarrow \tau) \rightarrow \tau)$, the amount introduced type variables would be $O(n)$.

Listing 3.2: Nested pairs with different type variables

```

1 fun dup a f = f a a;
2 fun one z = dup z;
3 fun two z = one (one z);
4 fun three z = two (two z);
5 ...

```

The trick to induce an exponential running time is with the pathological program in Listing 3.2. The type of **dup** is $a \rightarrow (a \rightarrow a \rightarrow \tau) \rightarrow \tau$. The type of **one** is $z_1 \rightarrow (z_1 \rightarrow z_1 \rightarrow \gamma_1) \rightarrow \gamma_1$.

The type of **two** is a bit more complicated; first solve the inner application by introducing the fresh variable τ_3 and let the right side be z_2 .

$$\text{unify}(z_1 \rightarrow (z_1 \rightarrow z_1 \rightarrow \gamma_1) \rightarrow \gamma_1, z_2 \rightarrow \tau_3) = \{z_2 \mapsto z_1, \tau_3 \mapsto (a \rightarrow a \rightarrow b) \rightarrow b\}$$

Furthermore solve the outer application with the fresh variable τ_4 .

$$\begin{aligned} & \text{unify}(z_1 \rightarrow (z_1 \rightarrow z_1 \rightarrow \gamma_1) \rightarrow \gamma_1, (a \rightarrow a \rightarrow b) \rightarrow b \rightarrow \tau_4) \\ &= \{z_1 \mapsto (a \rightarrow a \rightarrow b) \rightarrow b, \tau_4 \mapsto (z_1 \rightarrow z_1 \rightarrow \gamma_1) \rightarrow \gamma_1\} \end{aligned} \quad (3.1)$$

$$z_2 \rightarrow (((z_2 \rightarrow z_2 \rightarrow \gamma_3) \rightarrow \gamma_3) \rightarrow ((z_2 \rightarrow z_2 \rightarrow \gamma_4) \rightarrow \gamma_4) \rightarrow \gamma_2) \rightarrow \gamma_2).$$

Observe that the polymorphic function f kk

□

Bibliography

- [1] Alonzo Church. An unsolvable problem of elementary number theory. *American journal of mathematics*, 58(2):345–363, 1936.
- [2] Alonzo Church. *The calculi of lambda-conversion*. Number 6. Princeton University Press, 1985.
- [3] B Jack Copeland. The church-turing thesis. 1997.
- [4] Luis Damas. Type assignment in programming languages. 1984.
- [5] Luis Damas and Robin Milner. Principal type-schemes for functional programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 207–212, 1982.
- [6] BJ Heeren, Jurriaan Hage, and S Doaitse Swierstra. Generalizing hindley-milner type inference algorithms, 2002.
- [7] Roger Hindley. The principal type-scheme of an object in combinatory logic. *Transactions of the american mathematical society*, 146:29–60, 1969.
- [8] William A Howard. The formulae-as-types notion of construction. *To HB Curry: essays on combinatory logic, lambda calculus and formalism*, 44:479–490, 1980.
- [9] Thomas Johnsson. Lambda lifting: Transforming programs to recursive equations. In *Conference on Functional programming languages and computer architecture*, pages 190–203. Springer, 1985.
- [10] Robin Milner. A theory of type polymorphism in programming. *Journal of computer and system sciences*, 17(3):348–375, 1978.
- [11] Dana Scott. A system of functional abstraction, 1968. lectures delivered at university of california, berkeley. *Cal*, 63:1095, 1962.