

WYE: A FUNCTIONAL LANGUAGE

VERSION 1.0.0

ADITYA GOMATAM

July 10, 2024

Abstract

To learn more about functional languages and compilers, I decided to build a functional language. This language is intended not only to have a minimal set of features so that it be relatively easy to implement, but also to be useful. As functional programming is based upon the lambda calculus, I have decided to name the language Wye, as y is λ upside-down.

This document will compile the specification of Wye: its syntax, semantics, type system, intermediate representation, register allocation, and finally its translation into an assembly language of choice.

1 Introduction

Wye is a statically-typed functional language based on the lambda calculus. It utilizes a Hindley-Milner type system, which enables both parametric polymorphism and type inference in Wye. For clarity of your code, and to aid the compiler, however, you may annotate the type of Wye variables. Wye also allows the declaration of custom types.

1.1 Type Classes

The original ideas for Wye included the ability to create Haskell-like Type Classes, between one could establish subtype relations, and moreover, by which type variables could be bound (for example, one could annotate that a function argument generalize over all 'Ordered types, instead of all types). Due to the added complexity of implementing such a system, Type Classes are deferred to a later version of Wye.

1.2 List Comprehension

List comprehension should not be too hard to implement compared to the rest of this project, but it will require some additions to the grammar, and so in keeping with the ideal of minimal set of features, are also deferred to a future version.

1.3 User Input

In Wye version 1.0.0, users may only output to the screen. Taking user input is slightly more complicated and is thus deferred to a future version of Wye.

1.4 Type Aliases

Similar to TypeScript, it would be nice if the user could do something like the following to create an abbreviation for a type:

```
# just a tuple-type of length 2
type Pair 'a 'b = ('a, 'b);
```

2 Syntax

A Wye program is stored in a `.y` file and follows the below grammar. A Wye program is, in the most abstract view, simply a sequence of Wye statements.

Witten below in Backus-Naur form is roughly the Wye grammar. The start symbol is *Program*. Nonterminals start with a capital letter. Terminals are written in **this font**. *pat** denotes the Kleene star: zero or more repetitions of the pattern *pat*. *pat+* is a shorthand for *pat pat**. *pat?* denotes that *pat* may or may not occur. $(pat_1 \dots pat_n)$ groups pat_1, \dots, pat_n into a new pattern. $pat1|pat2$ means that exactly one of *pat1* or *pat2* may occur. Beware that `|` and `|` are different – the first is a Wye token. \langle These brackets \rangle are used to annotate certain grammar rules.

2.1 Grammar

<i>Program</i>	→	<i>Statement</i> *
<i>Statement</i>	→	<i>LetStatement</i> <i>TypeDeclaration</i>
<i>LetStatement</i>	→	let <i>Id</i> (: <i>Type</i> <i>Id</i> * ((<i>Id</i> : <i>Type</i> ->) * <i>Type</i>)) [?] := <i>Expr</i> ;
<i>TypeDeclaration</i>	→	type <i>TypeId</i> (' <i>Id</i>) * = (<i>TypeId</i> (with <i>Type</i>) [?]) (<i>TypeId</i> (with <i>Type</i>) [?]) * ;
<i>Expr</i>	→	<i>IntLiteral</i> <i>FloatLiteral</i> <i>StringLiteral</i> <i>List</i> <i>Tuple</i> <i>Id</i> <i>BuiltinOp</i> <variable/func within current scope> <i>TypeId</i> (with <i>Expr</i>) [?] <Type variant> <i>Expr Expr</i> <function application> <i>Expr BuiltinOp Expr</i> <reserved binary operator> match <i>Expr</i> { (<i>Pat</i> => <i>Expr</i> ;) * <i>Pat</i> => <i>Expr</i> (;?) } if <i>Expr</i> then <i>Expr</i> else <i>Expr</i> (<i>Expr</i>) { <i>Statement</i> * <i>Expr</i> } <let .. in block> print <i>Expr</i> error <i>Expr</i>
<i>Type</i>	→	int float string <i>TypeId</i> (<i>Type</i>) * [<i>Type</i>] ((<i>Type</i> ,) + <i>Type</i>) ' <i>Id</i> <type variable> <i>Type</i> -> <i>Type</i> <function type>
<i>List</i>	→	[(<i>Expr</i> ,) * <i>Expr</i>] []
<i>Tuple</i>	→	((<i>Expr</i> ,) + <i>Expr</i>)
<i>BuiltinOp</i>	→	+ - * / // :: < <= > >= == !=
<i>Pat</i>	→	<i>Expr</i> _ <i>Id</i> :: <i>Id</i> <head : tail list destructuring> <i>TypeId</i> (with <i>Id</i>) [?] [(<i>Pat</i> ,) * <i>Pat</i>] [] ((<i>Pat</i> ,) + <i>Pat</i>) ~ <i>Pat</i> <pattern negation>

The exact substitution rules for the *Literal* patterns listed on the first line of the *Expr* substitution rules is omitted to avoid boring the reader. If you are reading this, you know what a float is. Similarly, the rules for *Uppercase* and other character sets are omitted. Note that Wye uses ASCII as its character set for identifiers and type names, though strings may contain Unicode characters.

Ids and *TypeIds* should not be any of the builtin keywords such as `int`, `float`, `string`, `print`, `match`, `with`, `then`, and so on. *Ids* begin with a lowercase letter after as many underscores as you want, and *TypeIds* begin with an uppercase letter, also after as many underscores as you want. Thus, *Ids* and *TypeIds* cannot conflict with each other within their scope, and the parser does not find any ambiguity here.

Notice that `bool` is not a builtin type. This is because the *TypeDeclaration* system of Wye is used to define it in the Wye prelude. Following that, `if then else` is just syntactic sugar for `match` on the enum variants of `bool`.

One-line Wye comments begin with `#` and mark all following text until the next carriage return or newline as whitespace. Multiline Wye comments begin with `(*` and end with `*)`.

In Wye, the application of functions is always written in postfix notation, except for certain reserved binary operations (such as `+`) that may be written in infix notation. These binary operations are, under the hood, translated into postfix notation.

3 Semantics

3.1 Statements

Currently, there are two types of statements in Wye | `let` statements and type declarations. The following are valid examples of Wye statements:

3.1.1 Let Statements

Wye `let` statements declare constants (referred to throughout this document as variables for consistency with common notions of `let`) at the top-level, or within blocks. In Wye, functions are first-class – that is, they occupy the same status as “regular” types. They can be declared as variables, passed as function arguments, and so on.

Variables declared using `let` statements are visible throughout their entire scope, both backwards and forwards. Blocks create nested scopes. The variables declared in such nested scopes are not visible in outer scopes. Function arguments shadow over the outer scope. Here are some examples:

```
# variable (actually a constant)
let x = 4;
# type-annotated variable
```

```

let y: int = 4;
# function
let plus_4 x = (+ x 4);
# type-annotated functions
let plus_4 x: int -> int = (+ x 4);
let id: x: 'a -> 'a = x;
let into_tup x: int -> y: float -> z: string -> (int, float, string)
  = ((plus_4 x), y, (id z));
# blocks
let y = {
  5
};
let doublesum lst: [int] -> int = {
  let sum lst: [int] -> int = match lst {
    [] => 0;
    x :: xs => (+ x (sum xs));
  };
  (* 2 (sum lst)) # no semicolon
}; # semicolon

```

It would be slightly uncomfortable to express this within the grammar, but Wye reserves a `let` statement of the following form at the top-level to be its `Main` statement – that is, the expression on the right hand side is run when the `.y` file is executed from the command line:

```

let Main = print (plus_4 6);

```

Apart from `Main`, every function in Wye is pure, i.e., it will not affect the scope outside of the function. There are no global mutable variables, and the only function that is allowed to call other functions that take input or output to the terminal is `Main`.

3.1.2 Type Declarations

Type declarations in Wye are akin to enums in Rust. Variants of the type are separated by `|` and one may choose to give any variant an optional field. Unlike Rust, wherein enums can have multiple fields, Wye allows only one field. However, one can declare this field to have a tuple type, which will allow the variant to behave as if it holds multiple fields. Moreover, user-defined types can be *polymorphic* – that is, they may have a type argument. See the following examples:

```

type bool = false | true;
type Option 'a = None | Some with 'a;
type binary_tree 'a = Leaf
  | Node with ('a, binary_tree 'a, binary_tree 'a);

```

One can think of polymorphic user-defined types as “type functions” that take in a type parameter and output an instance of a type. In the above example `Some with 4` would take in the type `int` of the expression `4` and output an instance of `Option int`.

3.2 Expressions

The right hand side of every let statement in Wye must be an expression. Wye version 1.0.0 supports only a few expressions. Most of them exist in other languages and are fairly self-explanatory. This section will go over the less familiar ones.

3.2.1 Function Application

As a language that implements the lambda calculus, Wye supports function application, which expressed as the following grammar rule:

$$Expr \rightarrow Expr Expr$$

For example:

```
let f x: int -> y: int -> int = (+ x y);  
let plus_4 y: int -> int = f 4;
```

Functions in Wye are *curried*. That is, the definition of `func` actually looks something like this under the hood:

```
function f(int x) {  
  function g(int y) {  
    return x + y;  
  }  
  return g;  
}
```

Thus, in defining `plus_4` as `f 4`, `plus_4` is actually like the function `g`, which has captured the value of `4` from an outer scope, and which on input `y` returns `4 + y`.

This is also the reason why Wye's function types include an arrow between arguments. The function `f` has type `int -> int -> int`, the function `f 4` has type `int -> int`, and the value obtained from computing `plus_4 5` has type `int`.

3.2.2 List Construction

Just like Lisp, OCaml, Haskell, and many other functional languages, lists are an important part of Wye. As such, Wye provides list construction via the "cons" operator `::`. Wye also supports pattern matching of lists using cons notation, which is talked about in 3.2.5.

The cons operator operates as follows: if x is a value of type t and l is a list of type $[t]$, with value $[a_1, a_2, \dots, a_n]$ where $n \geq 0$, then $x :: l$ denotes the list $[x, a_1, \dots, a_n]$.

The list construction expression falls under the $Expr BuiltinOp Expr$ grammar substitution rule for $Expr$.

3.2.3 Using variants of a user-defined type

If t is a custom type declared via a *TypeDeclaration* statement, with variants t_1, \dots, t_n , then one may access the variant t_i as t_i . In version 1.0.0 of Wye, for simplicity of the parser, variants with the same name across different user-defined types is not allowed. Thus, the following Wye code compiles:

```
# synonyms of many
type SynMany = Numerous | Several | Plenty;
# synonyms of abundant
type SynAbundant = Copious | Profuse;
let x = Plenty;
let x = Copious;
```

but this does not:

```
# synonyms of many
type SynMany = Numerous | Several | Plenty;
# synonyms of abundant. ERROR - conflicting variants
type SynAbundant = Plenty | Copious | Profuse;
```

If the type variant has a field, then one should use **with** syntax to define the value of the field. For example:

```
type Option 'a = None | Some with 'a;
let x: Option int = Some with 4;
```

3.2.4 Match constructs

The **match** construct allows control flow of the program to be determined dynamically at runtime, based on whether or not a particular expression matches a pattern. The technical syntax of the **match** construct is as follows, with match branches separated by a semicolon (;):

$$\text{MatchExpression} \rightarrow \text{match Expr} \{ (Pat \Rightarrow Expr ;) * Pat \Rightarrow Expr (; ?) \}$$

or equivalently:

$$\text{MatchExpression} \rightarrow \text{match Expr} \{ Pat_1 \Rightarrow Expr_1 ; \dots Pat_n \Rightarrow Expr_n (; ?) \}$$

where $n \geq 1$.

The operation of the **match** construct is as follows: the *Expr* following **match** is first evaluated to a value v . It will then check at runtime whether v against each of Pat_1 to Pat_n in order. Let i be the smallest number from 1 to n such that v matches Pat_i (what it means to “match” a pattern will be discussed later). This branch will be taken, and only $Expr_i$ will be evaluated. If v matches *none* of the patterns Pat_1, \dots, Pat_n , Wye will throw a runtime error.

Each of $Expr_1, \dots, Expr_n$ must evaluate to a value of the same type under all possible executions. Thus, the **match** expression will evaluate to something of exactly one type, under all possible executions.

3.2.5 Patterns in a Match construct

The patterns that one can match against in a **match** are described by the following rules:

$$\begin{aligned}
 Pat \rightarrow & Expr \mid _ \\
 & \mid Id :: Id \\
 & \mid TypeId (with Id)? \\
 & \mid [(Pat ,) * Pat] \mid [] \\
 & \mid ((Pat ,) + Pat)
 \end{aligned}$$

Throughout the subsequent discussion, let v denote the value to which the $Expr$ after the **match** keyword evaluates. The following list describes when a particular pattern is matched:

- If $Pat \rightarrow Expr$, then the expression on the right is evaluated to some w , and v matches this expression iff there is exact equality of v and w . Note equality in Wye is always *by value*. This implies that if two expressions do not have the same type, they cannot ever be equal.
- In particular note that $Expr \rightarrow Id$ is possible. Let x denote the Id that occurs. If this x is already present in the outer scope of the **match** construct, then the value of this x is taken and exact equality is checked between the value of x and v . If such x does not exist, Wye should catch the undefined symbol at compile-time.
- If $Pat \rightarrow _$ then v always matches the pattern. $_$ is a wildcard in Wye.
- If $Pat \rightarrow Id :: Id$, then v matches the pattern iff v is a list of type $[t]$ with at least one element. This pattern type is better explained with the following example:

```

match e1 {
  x :: s => e2
}

```

Suppose the expression $e1$ evaluates to the value v . Then v matches $x :: s$ iff:

1. v is of some list type $[t]$
2. v has at least one element

If v matches the pattern, then the *new variable* x will be set to the first element of v (the “head” of v) and the *new variable* s will be set to the list containing the remaining elements of v , in the original order (that is, s will be set to the “tail” of v). x and s will be available for use in the scope of e_2 ONLY, and will shadow over variables with the same name in the outer scope. v will always fail to match such a “list-construction” pattern if it is a list of 0 elements. Note that if v is a list of 1 element, x will always be set to that element, and s will be an empty list.

- If $Pat \rightarrow TypeId$ then v must evaluate to a variant of a type defined in a *TypeDeclaration* statement. For example, the following pattern matching can be done:

```
type Operation = Add | Sub | Mul | Div;
let to_string op: Operation -> string = match op {
  Add => "add";
  Sub => "subtract";
  Mul => "multiply";
  Div => "divide"
}
let x = Operation.Mul;
let v = to_string x; # v is "multiply"
```

- Similarly, if $Pat \rightarrow TypeId \text{ with } Id$, then b must evaluate to a variant of a type defined in a *TypeDeclaration* statement *that has a field*. In this case, the semantics are similar to the list destructuring pattern $x :: s$. For example:

```
type Option 'a = None | Some with 'a;
let print_optional opt: Option string -> string = match opt {
  None => print "nothing!";
  # assume that cat concatenates two strings
  Some with x => print (cat "some with" x);
}
```

Thus, if x denotes the *Id* that occurs after the *with* keyword, then x is a new variable that is visible within the scope of the expression after the \Rightarrow . This *Id* will shadow over a variable with the same name from the outer scope.

- If $Pat \rightarrow []$ then v matches the pattern iff it is an empty list (and it type checks).
- If $Pat \rightarrow [(Pat,) * Pat]$, then v must be a list type in order to match the pattern. Let p denote this pattern. Then v matches p iff:
 1. v is a list of the same length as p (note in this case, $|p| \geq 1$)
 2. if v_i is the i^{th} element of v and p_i is the i^{th} pattern in p , then for each i , v_i must match p_i

- If $Pat \rightarrow ((Pat,) + Pat)$, then v must be a tuple type in order to match the pattern. Let p denote this pattern.. Then, v matches p iff:
 1. v is a tuple of the same length as p (note in this case, $|p| \geq 1$)
 2. if v_i is the i^{th} element of v and p_i is the i^{th} pattern in p , then for each i , v_i must match p_i (in type as well).
- If $Pat \rightarrow \sim Pat$, let p be the LHS Pat and q be the RHS Pat . Then v matches p iff v does not match q .

3.2.6 Function overloading

Wye does not support function overloading of user-defined functions. However, certain functions, such as `+` are defined on multiple types, which the compiler is able to handle. For an interested reader as to how one would handle function overloading well in a functional language, look up Haskell's "type classes" and "kinds".

3.2.7 Builtin Operations

There are 6 builtin operations in Wye:

- `+` adds two `ints` or `floats` or concatenates two `strings` or lists. Wye does not support addition of `int` and `float`. An explicit cast must occur (more on this later).
- `-` an expression `a - b` results in the subtraction of the value `b` from the value `a`. Subtraction is supported only on `ints` and `floats`.
- `*` multiplication of `ints` or `floats`.
- `/` an expression `a / b` results in the division of the value `a` by the value `b`. In Wye, division is supported *only for floats*. If `b` is the value 0, then Wye will throw a division-by-zero error at runtime.
- `//` an expression `a // b` results in the floor division of the value `a` by the value `b`. In Wye floor division is supported *only when a is a float and b is an int*. If `b` is the integer 0, then Wye will throw a division-by-zero error at runtime.
- `::` as discussed in 3.2.2, `::` is the list construction operator in Wye for constructing a list via specifying its head and tail.

Wye's builtin binary operations can be written both in postfix notation (for example, division can be written as `(/ a b)`), or in infix notation `(a / b)`.

3.3 Types

3.3.1 Basic types

Wye has essentially 3 builtin primitive types, with support for constructing list and tuple data structures. The three primitive types are:

int: a 64-bit signed integer.

float: a 64-bit floating point number.

string: essentially a list of 4-byte unicode characters.

The specification for the types of lists and tuples has already been demonstrated earlier.

3.3.2 Polymorphic types

Just like OCaml and Haskell, Wye's type system is based on the Hindley-Milner type system and thus supports polymorphic types. A polymorphic type is essentially a type of the form $\forall \alpha \sigma$ where σ is a type signature in which α is a free variable. Such type signatures have been referred to earlier, for example in the `Option` type:

```
type Option 'a = None | Some with 'a;
```

Or, in the identity function:

```
let id x: 'a -> 'a = x;
```

In Wye, the `'` token specifies that the *Id* immediately following it is a *type variable* which is implicitly quantified over in a type expression (an expression that evaluates to a type). In this way, the function `id` can take in any value of any type! The function's argument is polymorphic.

Whenever the `id` function is called, it actually implicitly takes in a *type argument* – namely, the type of its actual argument! Due to Wye's type system, this type is actually known at compile-time, and thus type checking can be adequately performed via substitution of the quantified type variable `a` with whatever type the function is called with. For example, the following code type-checks, with the type variable `a` being substituted with `string` and the type variable `b` being substituted with `int`:

```
type OptionalTuple 'a 'b = None | Some with ('a, 'b);
let f x: 'a -> y: 'b -> OptionalTuple 'a, 'b = Some with (x, y);
# j is of variant Some with (string, int)
let j = f "hello" 4;
```

3.3.3 Type inference

Wye is a statically typed language, and due to its Hindley-Milner type system, it is able to infer the types of expressions (and thus variables) at compile-time. The exact rules for the type system will be talked about in greater detail in the Type Checking section of this document. However, it is important to note that Wye expressions,

when not annotated with a type, are typed to the *most general possible type*. For example, in the following code:

```
let cons x lst = x :: lst;
```

the function `cons` type checks to the functional type `'a -> ['a] -> 'a`.

For more information, the interested reader should check out the [the Hindley-Milner type system](#).

4 Type Checking

In order to verify at compile-time that a program accomplishes the author's intention, it is helpful to type-check the program at compile-time. Type checking is the process of verifying that all operations in the program are performed on values of appropriate types. Type checking can be performed at compile-time only if the types of every Wye expression are known at compile-time. Thus, if the Wye compiler cannot figure out the type of a Wye expression, the compiler will refuse to compile the input program. Type checking is achieved via deductive reasoning on the input program based on certain deductive rules. In version 1.0.0 of Wye's compiler, type checking is implemented via a depth-first traversal of the Abstract Syntax Tree of an input Wye program.

4.1 Notation

The rules by which Wye performs type checking are written in the following form:

$$\frac{\text{hypothesis}_1, \dots, \text{hypothesis}_n}{\text{conclusion}}$$

In order to perform typing, a context is necessary. For example, if a function `f` is applied in an expression, then `f` must be defined somewhere. The set of variables, functions, and types available in scope for use in an expression is referred to as the *context* and notated as Γ .

The symbol \vdash means "entails" or "proves." We say $\Gamma \vdash e : \tau$ if the context Γ and the deductive rules of Wye's type system prove that the expression e has type τ . If interested in the origins of this symbol, the reader should look into soundness and completeness in mathematical logic.

A comma between hypotheses indicates a logical and. That is, $\text{hypothesis}_1, \text{hypothesis}_2$ is logically equivalent to $\text{hypothesis}_1 \wedge \text{hypothesis}_2$. Moreover, Γ can be considered a hypothesis itself, an assertion that all declared variables and functions have their specific types, along with assertions about the available types. Thus, the notation $\Gamma, \text{hypothesis} \vdash e : \tau$ may be used to express that $\Gamma \wedge \text{hypothesis}$ proves that $e : \tau$.

In general, τ will be used to indicate a monomorphic type (i.e., a "concrete" type like `int`), and σ will be used to indicate a polymorphic type. Any type with a type

variable, i.e., with an implicit quantifier over a type parameter, is a polymorphic type.

Also important to consider are the free type variables in a type expression. For example in the identity function:

```
let f x: 'a -> 'a = x;
```

within the expression on the right hand side of the $=$ (namely, x), the type of the parameter x is free. It is equal to some a which is only concretized when f is applied to some parameter.

On the other hand, there are no free variables in f 's type. f has the type $\forall a a \rightarrow a$ – all the type variables in this function definition have been fully quantified over.

The set of free variables in an expression or set of expressions S is notated as $\text{free}(S)$. free is defined as follows:

1. if α is a type variable, $\text{free}(\alpha) = \alpha$.
2. if C is a type function (a variant of a user-defined type with a field), then $\text{free}(C \tau) = \text{free}(\tau)$
3. if Γ is a context, then $\text{free}(\Gamma) = \bigcup_{x:\sigma \in \Gamma} \text{free}(\sigma)$ where the union is taken over all distinct polymorphic types of variables or functions in Γ .
4. for polymorphic types, $\text{free}(\forall \alpha \sigma) = \text{free}(\sigma) - \{\alpha\}$
5. for contexts, $\text{free}(\Gamma \vdash e : \sigma) = \text{free}(\sigma) - \text{free}(\Gamma)$

4.2 Type order

Before we give the rules for type inference, it is important to define an ordering on types. This allows Wye to infer the most general type possible for an expression (the “least” potential type in this type ordering), which creates a highly flexible yet still strict type system.

4.3 Typing expressions

5 Acknowledgements

The syntax and semantics of Wye are based on Rust, Haskell, and OCaml. I thank the developers of these languages for making it much easier for me to learn about and write a compiler. I thank Jens Palsberg of UCLA for teaching me the fundamentals of compiler implementation.