



DUBLIN INSTITUTE
of TECHNOLOGY
Institiúid Teicneolaíochta Bhaile Átha Cliath

IMPLMENTING A INTERPRETER FOR A SCRIPTING LANAGUGE USING HASKELL

FINAL YEAR PROJECT REPORT

Zhen LAO

`zhen.lao@student.dit.ie`

Supervisor: Richard LAWLOR

2nd Reader: Cindy LIU

April 1, 2011

This Report is submitted in partial fulfillment of the requirements for the
award of the degree of **BSc Computer Science** of the School of
Computing, College of Sciences and Health, Dublin Institute of Technology.

Abstract

In this thesis,

Keywords:programming language YUN

Declaration

I **Zhen Lao** hereby declare that the work described in this dissertation is, except where otherwise stated, entirely my own work and has not been submitted as an exercise for a degree at this or any other university.

Signed _____
Zhen Lao

Acknowledgements

I would like to thank my supervisor Richard Lawor, for his valuable advice and useful suggestions on my project.

I am also deeply indebted to all the other tutors and teachers in Computer Science for their direct and indirect help to me.

Special thanks should go to my friends who have put considerable time and effort into their comments on the draft.

Contents

1	Introduction	6
1.1	Objective and Motivation	6
1.2	Benefits of using Haskell	6
1.3	Development Methodology	7
2	Compiler and Interpreter Technologies	8
2.1	Parsing technologies	8
2.1.1	Formal Grammar	8
2.1.2	The Hierarchy of Grammars	9
2.1.3	Backus–Naur Form and Extended Backus–Naur Form	9
2.2	Parser Generator Haskell Happy	10
2.3	Monadic Parsing using Parsec	10
2.4	Lexical analysis	10
2.4.1	Regular Expression and token	11
2.4.2	The Lexer Generator Alex	11
3	Monads and Haskell	12
3.1	Haskell and Category Theory	12
3.2	Monadic Function	13
3.2.1	Data Type Constructor	13
3.2.2	Monadic Function	13
3.3	Monads	13
3.3.1	IO Monad	14
3.3.2	State Monad	15
3.3.3	List Monad	15
3.4	Using Monad Operator to Combine Monadic Function	15
3.5	Type system in Haskell	15
3.6	Do Notation - Reinventing a imperatilanguage language	15

4	Monad Transformer	16
4.1	State Transformer	16
4.2	Error Transformer	16
4.3	Putting All Together	16
5	Language Interpreter Design	17
5.1	Imperative and Declarative language	17
5.2	Type System	17
5.2.1	Dynamic Typing and Static Typing	17
5.2.2	Strong Typing and Weak Typing	17
5.3	Problem and resolution in writing BNF/EBNF rules	18
5.3.1	Shift//Reduce Problem	18
5.3.2	Reduce//Reduce Problem	19
5.4	Syntax Design	19
5.4.1	The Main Structure	19
5.4.2	Statements	20
5.4.3	Generic Expression	22
5.4.4	Code Block	25
6	Language Implementation	27
6.1	Data Type	27
6.1.1	Tokenizer Type	27
6.1.2	Parser internal data type	28
6.1.3	Interpreter internal data type	28
6.2	Symbol Table and Parse Tree	28
6.3	Generic Expression Interpreter	28
6.4	Statement Interpreter	28
6.4.1	Function Invocation	28

List of Figures

4.1	Diagram illustrating a state monad	16
5.1	Module Syntax Diagram	19
5.2	Main Function Syntax Diagram	19
5.3	Function Syntax Diagram	20
5.4	Statement Syntax Diagram	21
5.5	Assign Statement Syntax Diagram	22
5.6	Expression Syntax Diagram	23
5.7	Expression Syntax Diagram	24
5.8	List Member Syntax Diagram	24
5.9	While Block Syntax Diagram	25
5.10	If Block Syntax Diagram	25
5.11	for Block Syntax Diagram	25

Chapter 1

Introduction

1.1 Objective and Motivation

The objective of this project is to develop a weak-type interpreted language using Haskell. This language is able to support the following feature,

- basic for loop and while loop
- basic if-else statement
- functional invocation
- arbitrary dimension list
- polymorphic list

Furthermore, in project, the monadic design approach is applied as Haskell is different from other object oriented language.

1.2 Benefits of using Haskell

Haskell is an advanced purely-functional programming language. By applying the use of Haskell to this project, I have significantly reduced the coding time and spent most of my time to the design phase.

A pure function is a function that accepts an input and generates an output. In Object-Oriented language, programs are constructed using classes and instances which encapsulate computation and state. Haskell programs are constructed by functions as functions are the first class members in Haskell. Typically the main function is defined in terms of other functions, which in turn are defined in

terms of still more functions, until at the bottom level the functions are language primitives. All of these functions are much like ordinary mathematical functions. [Why Functional Programming Matters] [Functional Programming with Overloading and Higher-Order Polymorphism]

1.3 Development Methodology

Agile development methodology is used in the entire development process. This project has been initially identified to multiple iterations and each iteration contains three major stages including research ,development and testing.

Chapter 2

Compiler and Interpreter Technologies

2.1 Parsing technologies

2.1.1 Formal Grammar

Mathematically, formal grammar consists of:

- a finite set of terminal symbols.
- a finite set of non-terminal symbols.
- a finite set of project rules.
- a start symbol.

[?] [Three Models for the Description of Language] From the formal grammar definition, legitimate production rules can be written as

$$S \mapsto aS \text{ and } S \mapsto ab$$

In this example, we can assume that the grammar consists of two projection rules and the starting symbol is S . The terminal symbols are lower letters $\{a, b\}$. From this example, If we start from the either rule 1 or rule 2, we could derive a grammar of $\{a^n b | n > 1\}$, which can be enumerate like $\{aab, aaab, aaaab, \dots\}$.

In addition, we are able to write all the production rule from the given abstract language, the language like

2.1.2 The Hierarchy of Grammars

Noam Chomsky has describe three model of grammar [”Three models for the description of language”] and this grammar model has significantly effect the design of computer programming language.

Chomsky define a set of rule upon the formal grammar and categorize them into different levels.

The Chomsky hierarchy consists of the 4 levels:

- Type-0 grammars (unrestricted grammar). It is a unrestricted grammars that include all possible grammar that are possible to recognize by Turning machine.
- Type-1 grammars (context-sensitive grammar).if all rules are of the form $\alpha A \beta \rightarrow \alpha \gamma \beta$ where $\alpha \beta \gamma$ are terminal symbols and A is non-terminal symbol.
- Type-2 grammars (context-free grammar).
- Type-3 grammars (regular grammar).

2.1.3 Backus–Naur Form and Extended Backus–Naur Form

The Backus-Naur Form(BNF) is a metalanguage to write the production rule that expressing the type-2 grammar (context-free grammar).It restricts the appearance of terminal and non-terminal in each side of the production equation.A canonical BNF production rule may like follow,

$$< symbol > ::= _expression_$$

The left side of the equation can only be non-terminal thus enclosed with $<>$.The right hand side can be terminals and non-terminals,a vertical bar ‘—’ is used to represent choice between terminal and non-terminals.

The Extended Backus–Naur Form (EBNF) and extension upon the BNF.Three regular expression qualifier is added to simplified some expression,they are,

- ? : which means that the symbol (or group of symbols in parenthesis) to the left of the operator is optional (it can appear zero or one times)
- * : which means that something can be repeated any number of times (and possibly be skipped altogether)

- $+$: which means that something can appear one or more times

[?]

Recursive rules of BNF like

$$1. \langle exp \rangle := \langle exp \rangle | sub$$

$$2. \langle exp \rangle := sub$$

that expressing a sequence of a particular syntactic element can be simplified using quantifier in EBNF as $\langle exp \rangle := sub^+$

2.2 Parser Generator Haskell Happy

Happy is a parser generator system for Haskell, similar to the tool ‘yacc’ for C. Like ‘yacc’, it takes a file containing an annotated BNF specification of a grammar and produces a Haskell module containing a parser for the grammar. [The Parser Generator for Haskell]

By using its own EBNF like syntax, used could write a parser description. The happy parser generator are able to recognize and compile it into Haskell source code.

2.3 Monadic Parsing using Parsec

In the early stage of this project, parse is build using parse C, Parsec is an industrial strength, monadic parser combinator library for Haskell. It can parse context-sensitive, infinite look-ahead grammars but it performs best on predictive (LL[Compilers: principles, techniques and tools.]) grammars. Combinator parsing is well known in the literature and offers several advantages to YACC or event-based parsing. [Parsec, a fast combinator parser]

Compared to parser generator, monadic parsing has two major benefits
 1. No need to learn additional parser generator grammar since parser combinator is written in the same language.
 2. parser can be adjust easily.

2.4 Lexical analysis

Before parsing, the lexical analyzer will scan the source code and generate a sequence of tokens.

2.4.1 Regular Expression and token

Tokens are defined by using regular expression,

2.4.2 The Lexer Generator Alex

In this project, the Alex Haskell Lexer generator is apply in generating token streams.

each token can be defined using regular expression.

Chapter 3

Monads and Haskell

3.1 Haskell and Category Theory

Category theory is a general theory that examine and organize mathematical object like set ,function,function domains Cartesian-set.

A Category C in category theory is defined below :

1. a collection of objects
2. a collection of arrows (often call morphism)
3. operations assigning to each arrow f an object $dom f$,its domain ,and an object $cod f$,its co domain.
4. a composition operator assigning to each pair of arrows f and g ,with $cod f = dom g$,a composite arrow $g \circ f : dom f \rightarrow cod g$, satisfying the following associative law:
For any arrow $f : A \rightarrow B, g : B \rightarrow C$, and $h : C \rightarrow D$ (with A,B,C and D not necessarily distinct),

$$h \circ (g \circ f) = (h \circ g) \circ f$$

5. for each object A, an identify arrow $id_a : A \rightarrow A$ satisfying the following identity law:
For any arrow $f : A \rightarrow B$,

$$id_a \circ f = f \text{ and } f \circ id_a = f.$$

[?]

Functions are the first member of the program in functional programming,since no size affect is not allow ,there should be a way to combine the

all kinds of functions to form a new function instead of just simply chain the input output of each function as the former will generate intermediate output.

For instance ,counting the file of java source code in current directory can be written as follow:

$$ls -al | grep * .txt | wc -l$$

To substantiate the this concept , let's use the map/fold fusion technique of Haskell as an example.

If we want to calculate the sum of the square of each element of a list eg. [1,3,4,6,7,9],the result of it is $1^2 + 3^2 + 4^2 + 6^2 + 7^2 + 9^2 = 192$.In Haskell ,we could use map and fold to address problem.

To avoid generating intermediate output from the first function to second function, the could rewrite the hold function using a single fold

The all map/fusion is is equivalent to $foldr f e . map g = foldr (\lambda xy -> f(gx)y) e$
therefore, the
 $sum_of_square = foldr (\lambda xy -> x^2 + y) 0$

3.2 Monadic Function

3.2.1 Data Type Constructor

3.2.2 Monadic Function

A monadic function is function that produce , however,monadic function like **putStr :: String -> IO ()** can not be combined using $(.) :: (b -> c) -> (a -> b) -> a -> c$.Monadic class constructor has tag

3.3 Monads

In Haskell,monad is used an abstract data type constructor to represent multiple kinds of computation such as a computation that will do IO action,or a computation that has state.Those computations are in-pure because that manipulate the outside world.In Haskell.Mathematically, monads are governed by set of laws that should hold for the monadic operations [A Gentle Introduction to Haskell, Version 98]. There are two basic law in monads ,they are bind return .The Monad class is defined as follow:

```

class Monad m where
  (>>=) :: m a -> (a -> m b) -> m b
  (>>)  :: m a -> m b -> m b
  return :: a -> m a
  fail   :: String -> m a

```

The return function can inject a value into monadic type. The bind function can combine two monadic function, one should be of type **m a** and another should be of type **a -> m b**.

Beside this two function, Haskell also provide other monadic operator which all derive from **return** and **bind**, they are:

```

liftM :: (Monad m) => (a1 -> r) -> m a1 -> m r
liftM2 :: (Monad m) => (a1 -> a2 -> r) -> m a1 -> m a2 -> m r
ap :: (Monad m) => m (a -> b) -> m a -> m b
(<=<) :: (Monad m) => (a -> m b) -> m a -> m b
\ $ :: (m a -> m b) -> m a -> m b

```

These monadic operation is define using the bind and return. For example, liftM is defined by bind and return like

```
liftM f m1 = do { x1 <- m1; return (f x1) }
```

Therefore, when defining a monad, only bind and return need to be specified.

Monadic Characteristic

For all monad instance, beside define three monadic operator, they must apply three compulsory monad laws:

```

"Left identity": return a >>= f  ≡  f a
"Right identity": m >>= return  ≡  m
"Associativity": (m >>= f) >>= g  ≡  m >>= (\x -> f x >>= g)

```

In monad instance, these monad laws will become a restriction of the operation when combining monadic function using monadic operation, these restriction will be discussed in following sections.

3.3.1 IO Monad

Haskell use IO monad to limit the IO sequence. Monadic operation are used to represent IO processing pipeline, One significant different between IO monad an other monad is that it does not provide escape function like

```
IO String ->String
```


3.3.2 State Monad

The State Monad is defined as follow

```
newtype State s a = State {runState :: s -> (a, s)}
```

3.3.3 List Monad

Haskell try to use list monad to represent a calculation that return multiple result.

3.4 Using Monad Operator to Combine Monadic Function

3.5 Type system in Haskell

3.6 Do Notation - Reinventing a imperatilan- guageve language

Chapter 4

Monad Transformer

Monad transformers offer an additional benefit to monadic programming: by providing a library of different monads and types and functions for combining these monads, it is possible to create custom monads simply by composing the necessary monad transformers. [Monad Transformers Step by Step]

4.1 State Transformer

A value of type `(ST a s)` is a computation which transforms a state index by type `s`, and delivers a value of type `a`. You can think of it as a box, like this

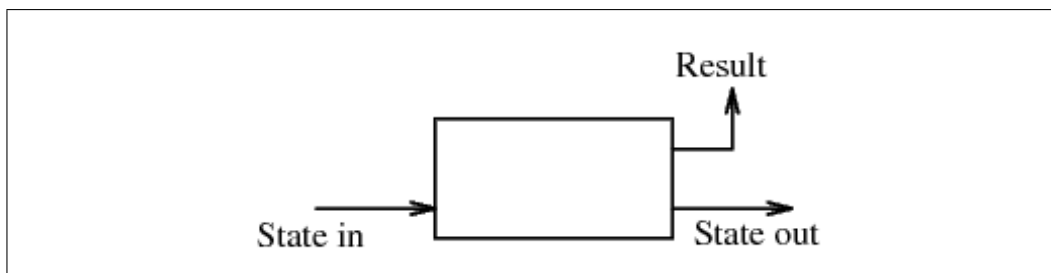


Figure 4.1: Diagram illustrating a state monad

[lazy functional state thread]

4.2 Error Transformer

4.3 Putting All Together

Chapter 5

Language Interpreter Design

5.1 Imperative and Declarative language

5.2 Type System

5.2.1 Dynamic Typing and Static Typing

A programming language is said to use static typing when type checking is performed during compile-time as opposed to run-time. For example , you have to specify the type explicitly and the compile will check the type correctness of a variable. A variable of a specified type can not assign to another value of other type.

Static Typing	Dynamic Typing
<pre>int a =1; /*a is of type int */ int a="a string"; /* its not valid to assign a string to a variable that has type int */</pre>	<pre>a=1; /* does not need to specified a type for this variable */ a ="a string"; /* it valid to change the type of the variable */</pre>

In my project, I used the dynamic typing scheme, which is , does not need to specify any type of variable and able to assign any type of primitives to a variable.

5.2.2 Strong Typing and Weak Typing

A language is said to be strong typing is that it place restriction in operation where data type can not be intermix.

Strong Typing	Weak Typing
<pre> a=123; /* a is a number */ b="123" /* b is a string */ c=a+b /* return type error */ </pre>	<pre> a = 123; b = "123"; c = a+b; /* either a will be convert to a string or b will be convert to a number */ </pre>

In my project,I have implemented a weak typing system .I have design an statement call generic expression , which allow different kinds of value to intermix with each other. An expression like "12343" + 1232 -324 can be parse as follow syntax tree.

5.3 Problem and resolution in writing BN-F/EBNF rules

For a parser , there are four operation it will do when encounters a terminal/non-terminal symbol.They are,

- Shift - push token onto stack
- Reduce - remove handle from stack and push on corresponding nonterminal
- Accept - recognize sentence when stack contains only the distinguished symbol and input is empty
- Error - happens when none of the above is possible; means original input was not a sentence

Conflicts arise from ambiguities in the grammar when two or more operations and rules that apply to the same sequence of input.[A final solution to the Dangling else of ALGOL 60 and related languages]

5.3.1 Shift//Reduce Problem

A shift conflict occurs if there are two or more rules that apply to the same sequence of input for the same operation reduce. This usually indicates a serious error in the grammar.

5.3.2 Reduce//Reduce Problem

A reduce/reduce conflict occurs if there are two or more rules that apply to the same sequence of input for the same operation reduce. This usually indicates a serious error in the grammar.

5.4 Syntax Design

5.4.1 The Main Structure

A module is the minimum executable unit in **yun**. It is composed by one main function and several functions. The main function is a entry point, it may invoke other functions.

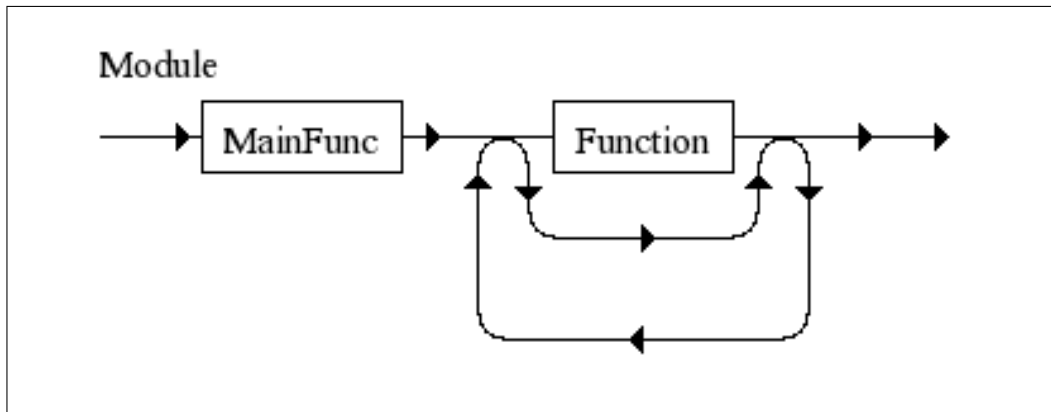


Figure 5.1: Module Syntax Diagram

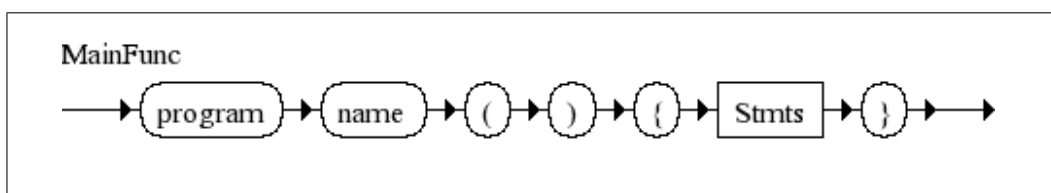


Figure 5.2: Main Function Syntax Diagram

Possible program code may look like follow,

```
1 program main ()
2 {
3     result = sum(1,2,3);
4     // more code
5 }
6 }
```

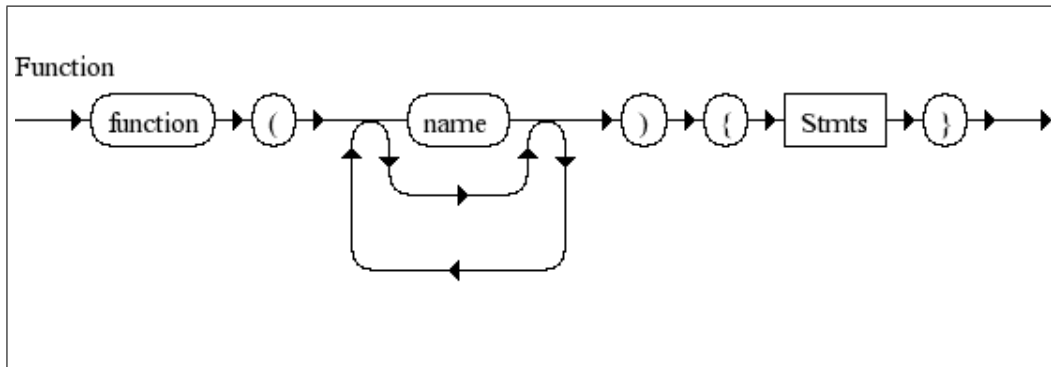


Figure 5.3: Function Syntax Diagram

```

7
8
9 function sum (var1 , var2 , var3)
10 {
11     // body of the functions
12 }
13
14 // may be more functions
  
```

5.4.2 Statements

Statements comprise the body of a function. A statement may be an assignment, break and continue sentence, return sentence. WhileBlocks, IfBlocks, ForBlocks are all statements. What's more, WhileBlock, IfBlock, ForBlock are recursively defined by statements as well. This allows nested for loops, nested while loops and etc.

The following code can be considered as a statement of the language

```

1 while ( )
2 {
3 } /* the while statement */
4
5 if ( )
6 {
7 }
8 else
9 {
10 } // the if else statement
11
12 for ( )
  
```

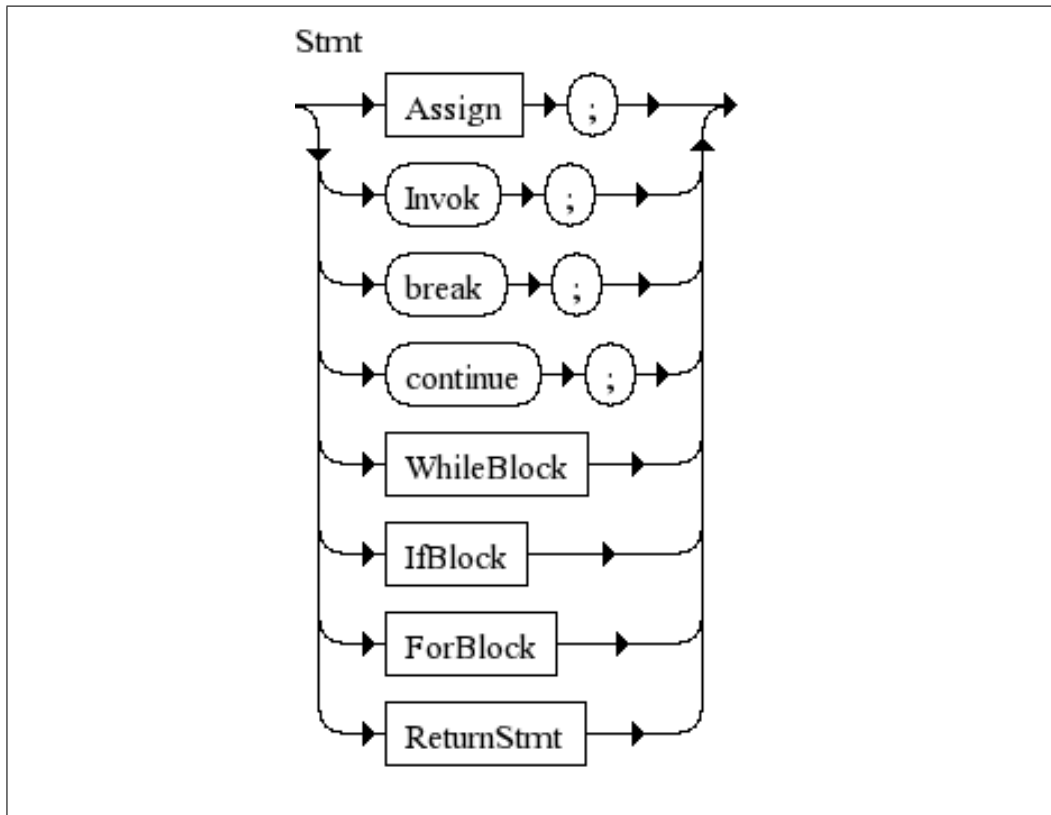


Figure 5.4: Statement Syntax Diagram

```

13 {
14     for ( )
15     {
16     }
17 } // nested for loop
18
19
20 a = 1 ; // assignment statement, assign a value or and an
           expression to a variable
21
22 fib(5); // function invocation statement
  
```

Assignment

There are three type of assignment in the language .They are ,

- Function invocation assignment.Assign the return result to a variable.

- Generic expression assignment. Assign the result of an expression to a variable.
- List assignment. Assign a list to a variable.

As **yun** is an imperative language, the right side of the assignment operator will be evaluation immediately. In other word, the language is strict.

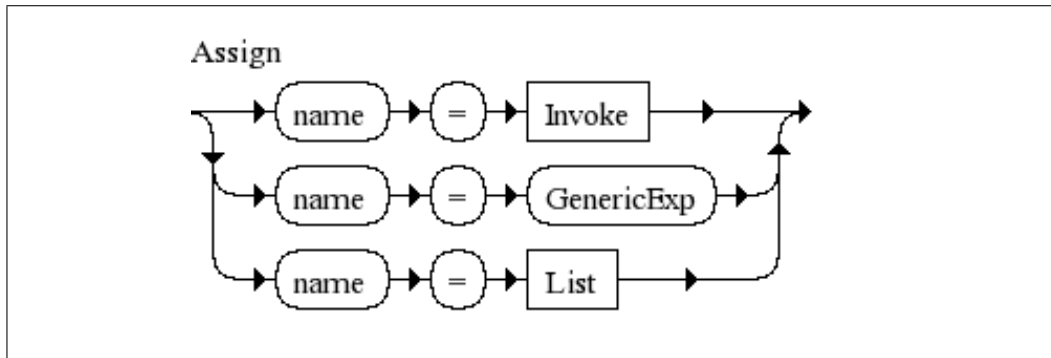


Figure 5.5: Assign Statement Syntax Diagram

5.4.3 Generic Expression

Generic expression represent the computation between primitives. As mention above, **yun** is a weak typing and dynamic typing language, the generic expression accepts variant types of primitive and will do the conversion internally.

Valid expression could be ,

```

1 "string"
2 1
3 -10
4 1.1
5 True
6 False
7 True
8 False // primitives
9
10 1 + "a string"
11 1 * 2 + 1
12 a && b
13 a || b
14 ! a
15 a[2] //get the third element of a list

```



```

16 a[1,2] // get the element in a multi dimension list
17
18 (a+b)+c // operation between variables
19 a[num]

```

List

The language support polymorphic list, the element of a list can be an expression or an list, thus the list support multi dimension list. If the element is an expression, the interpreter will evaluate it and store its result value in its internal format.

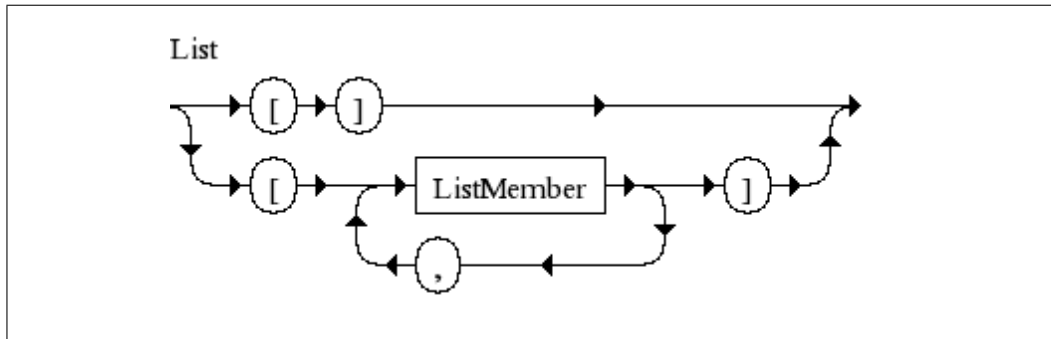


Figure 5.7: Expression Syntax Diagram

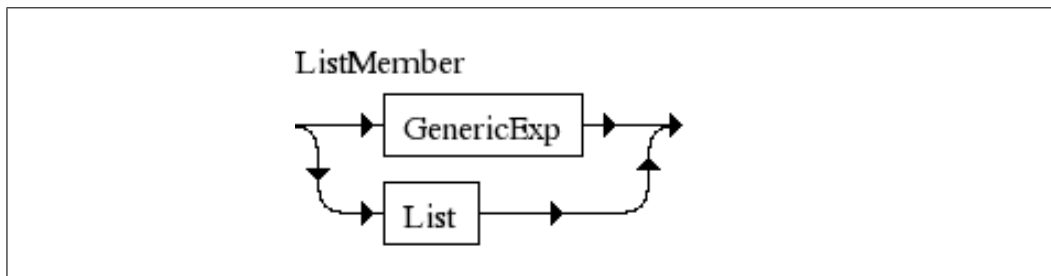


Figure 5.8: List Member Syntax Diagram

Code Example

```

1 a = [1,2,3];
2 a = ["string",1,2,3];
3 a = [var1,var2,"other",var1+var2,True&&False];
4 a = [[1,2,3,4,5,6,7] , 1, 23.32, 5];
5 a = [1,2,3,4,5],[1,2,3,4],[3,45,43]]; // declare a multi
    dimension list
6
7 member = a[1] // get the second element of a list

```

5.4.4 Code Block

Code blocks including if block, while block and for block are statements. Code block may contain other statements which allow nested code block to be defined.

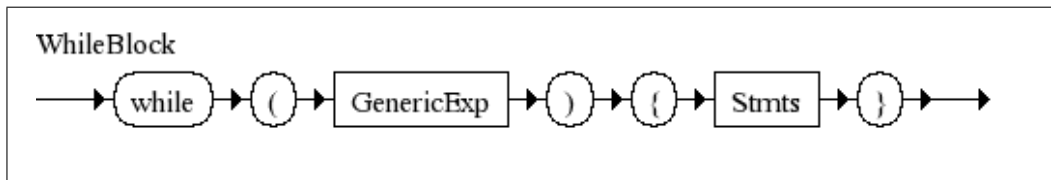


Figure 5.9: While Block Syntax Diagram

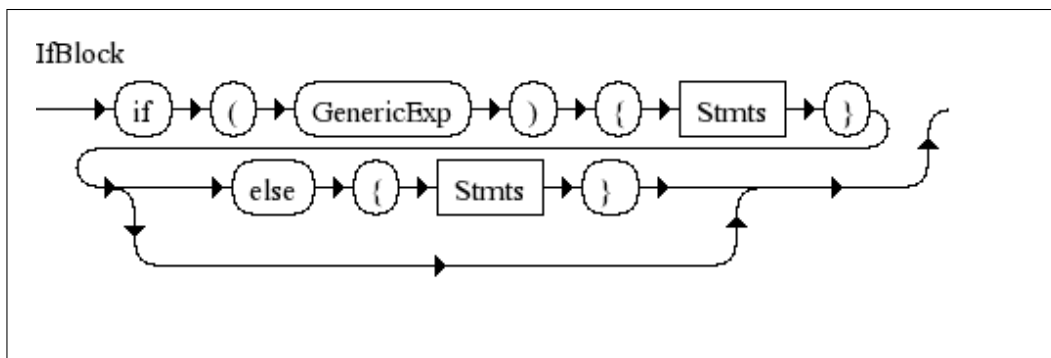


Figure 5.10: If Block Syntax Diagram

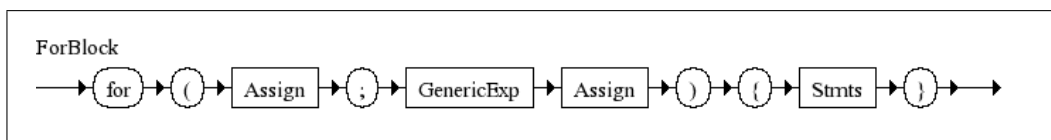


Figure 5.11: for Block Syntax Diagram

Code Example,

```
1 if (a < 1){  
2     a = a+1;  
3     if ( a == 1){  
4         a = a-1;  
5     }  
6 }  
7 else{  
8     a = a-1;  
9 } // nested if loop  
10
```

```

11
12 while (a<10){
13     a= a+1;
14 } // while block
15
16
17 for (i =0;i <10;i=i+1){
18     for (j=0;j <10;j=j+1){
19         sum=j+i ;
20     }
21 }//nested for block

```

Chapter 6

Language Implementation

6.1 Data Type

6.1.1 Tokenizer Type

All tokens are define using the Haskell Lexer Generator Alex. Tokens will be recognized and converted to Haskell data types. For example , a valid variable is comprised with one or more alphabetic characters and digit and the first character must be an alphabetic characters. This rule can be defined as ,

```
1 $digit = 0-9                — digits
2 $lowerCase = [a-z]
3 $alpha = [a-zA-Z]           — alphabetic characters
4
5 tokens :-
6   $alpha [$alpha $digit \_ \' \.]* { \s -> TName s }
7 -- data definition
8 data Token =
9     TName String |
10    TBool Bool |
11    TInt Int |
12    ... -- more definition
```

The variable **var1** will be parse into Haskell data type **TName "var1"**

By using the lexer , all the code will be generated into tokens streams.

```
1 program myProgram ()
2 {
3     var1 = 1;
4     result = var1 + "a string";
5 }
```

The following codes show the invocation of the lexer and the returned tokens stream.

```
1 lexer "program myProgram() { var1 = 1; result = var1 +  
    False; return 0; }"  
2 [TProgram,TName "myProgram",TOPB,TCPB,TOCB,TName  
    "var1",TAssign,TInt 1,TSC,TName "result",TAssign,TName  
    "var1",TPlus,TBool False,TSC,TReturn,TInt 0,TSC,TOCB]
```

6.1.2 Parser internal data type

The Parser internal data type represents a parse tree.

6.1.3 Interpreter internal data type

6.2 Symbol Table and Parse Tree

6.3 Generic Expression Interpreter

6.4 Statement Interpreter

6.4.1 Function Invocation