

---

# Compiler Theory and Practice

Coursework

Juan Scerri  
123456A

April 4, 2024

*A coursework submitted in fulfilment of study unit CPS2000.*



## Contents

Contents	i
Listings	iii
Report	1
§ 1    Lexer	1
1.1   Design & Implementation	1
§ 2    The Parser	10
2.1   Modified EBNF	11

2.2	Parsing & The Abstract Syntax Tree . . . . .	17
2.3	The Actual Parser . . . . .	18
2.4	Pretty? Printing . . . . .	27
§ 3	Semantic Analysis . . . . .	31
3.1	Phases & Design . . . . .	31
3.2	The Environment Tree . . . . .	32
3.3	Semantic Analysis . . . . .	40
3.4	The Symbol Type . . . . .	40
3.5	Type Checking . . . . .	43
§ 4	Code Generation . . . . .	50
4.1	Passing on the Environment Tree . . . . .	50
4.2	Reordering . . . . .	51
4.3	Function Declarations . . . . .	52
4.4	Scopes & Variable Declarations . . . . .	52
4.5	Expressions . . . . .	52
4.6	Branching . . . . .	52
§ 5	Attributions . . . . .	52

## Listings

1	DFSA Class Declaration (lexer/DFSA.hpp) . . . . .	1
2	Code specification of the comments in the LexerDirector (lexer/LexerDirector.cpp) . . . . .	7
3	Registration of the hexadecimal category checker (lexer/LexerDirector.cpp) . . . . .	7
4	Constructions of the Lexer (lexer/LexerBuilder.cpp) . . . . .	8
5	The updateLocationState() lexer method (lexer/Lexer.cpp) . . . . .	8
6	Error handling mechanism in the nextToken() lexer method (lexer/Lexer.cpp) . . . . .	9
7	The Runner constructor passes mLexer into the Parser constructor (runner/Runner.cpp) . . . . .	10
8	Mechanism for internally storing types within the compiler (parl/Core.hpp) . . . . .	14
9	The Primitive structure (parl/Core.hpp) . . . . .	15
10	An size unbounded type in C++ (parl/Core.hpp) . . . . .	15
11	The FunctionCall AST node class (parl/AST.hpp) . . . . .	17
12	The Binary AST node class (parl/AST.hpp) . . . . .	17
13	The Unary AST node class (parl/AST.hpp) . . . . .	17
14	The moveWindow() and nextToken() Parser methods (parser/Parser.cpp) . . . . .	19
15	The consume() Parser methods (parser/Parser.hpp) . . . . .	20
16	The Usage of the consume() method in the formalParam() generator method (parser/Parser.cpp) . . . . .	20
17	The peek() Parser method (parser/Parser.cpp) . . . . .	21
18	The abort functionality present in the codebase (parl/Core.hpp) . . . . .	21
19	The ParseSync exception and the error() method which kickstarts the synchronization process (parser/Parser.hpp) . . . . .	22
20	The synchronize() method in the Parser class (parser/Parser.cpp) . . . . .	23
21	The ifStmt() node generator method in the Parser class (parser/Parser.cpp) . . . . .	25
22	The main body of methods in the Parser class (parser/Parser.hpp) . . . . .	26

23	The parse segment of the <code>run()</code> method in the <code>Runner</code> class ( <code>runner/Runner.cpp</code> ) . . . . .	27
24	The <code>debugParsing()</code> method in the <code>Runner</code> class ( <code>runner/Runner.cpp</code> )	29
25	A segment of the pure virtual <code>Visitor</code> class ( <code>parl/Visitor.hpp</code> ) . . .	30
26	The <code>visit(core::PadRead*)</code> method in the <code>PrinterVisitor</code> ( <code>parser/PrinterVisitor.cpp</code> ) . . . . .	31
27	The <code>Environment</code> class with <code>mChildren</code> highlighted ( <code>backend/Environment.hpp</code> ) . . . . .	35
28	The <code>pushEnv()</code> method in the <code>EnvStack</code> class ( <code>analysis/EnvStack.cpp</code> )	37
29	The <code>popEnv()</code> method in the <code>EnvStack</code> class ( <code>analysis/EnvStack.cpp</code> )	37
30	The <code>pushEnv()</code> method in the <code>RefStack</code> class ( <code>ir_gen/RefStack.cpp</code> )	37
31	The <code>popEnv()</code> method in the <code>RefStack</code> class ( <code>ir_gen/RefStack.cpp</code> )	38
32	Definitions of <code>VariableSymbol</code> and <code>FunctionSymbol</code> ( <code>backend/Symbol.hpp</code> ) . . . . .	40
33	The <code>visit(FormalParam *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	41
34	The <code>visit(Variable *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	43
35	A part of the <code>visit(Binary *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	43
36	Another part of the <code>visit(Binary *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	44
37	Checking whether return statement is inside a function declaration in the <code>visit(ReturnStmt *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	45
38	Checking whether the return expression has the same type as the function return type in the <code>visit(ReturnStmt *)</code> method in the <code>AnalysisVisitor</code> class ( <code>analysis/AnalysisVisitor.cpp</code> ) . . . . .	45
39	The <code>visit(IfStmt *)</code> method in the <code>ReturnVisitor</code> class ( <code>analysis/ReturnVisitor.cpp</code> ) . . . . .	47
40	The <code>visit(Block *)</code> method in the <code>ReturnVisitor</code> class ( <code>analysis/ReturnVisitor.cpp</code> ) . . . . .	47
41	Getting the global environment from the <code>AnalysisVisitor</code> and passing it on for reordering and code generation ( <code>runner/Runner.cpp</code> ) .	50
42	The <code>visit(FunctionCall *)</code> method in the <code>TypeVisitor</code> class ( <code>ir_gen/TypeVisitor.cpp</code> ) . . . . .	51



# Report

## 1 | Lexer

### 1.1 | Design & Implementation

The lexer was split into three-main components. A DFSA class, a generic table-driven lexer, and a lexer builder.

#### The DFSA

The DFSA class is an almost-faithful implementation of the formal concept of a DFSA. Listing 1, outlines the behaviour of the DFSA. Additionally it contains a number of helper functions which facilitate getting the initial state and checking whether a state or a transition category is valid. These helpers specifically, `getInitialState()` is present since after building the DFSA there is no guarantee the initial state used by the user will be the same.

---

```
class Dfsa {
public:
    Dfsa(
        size_t noOfStates,
        size_t noOfCategories,
        std::vector<std::vector<int>>> const&
            transitionTable,
        int initialState,
        std::unordered_set<int> const& finalStates
    );

    [[nodiscard]] int getInitialState() const;

    [[nodiscard]] bool isValidState(int state) const;
```

```
[[nodiscard]] bool isValidCategory(int category) const;

[[nodiscard]] bool isFinalState(int state) const;

[[nodiscard]] int getTransition(
    int state,
    std::vector<int> const& categories
) const;

private:
    const size_t mNoOfStates;           // Q
    const size_t mNoOfCategories;       // Sigma
    const std::vector<std::vector<int>>
        mTransitionTable;               // delta
    const int mInitialState;             // q_0
    const std::unordered_set<int> mFinalStates; // F
};
```

---

Listing 1: DFSA Class Declaration (lexer/DFSA.hpp)

The only significant difference is the `getTransition()` functions. In fact, it accepts a vector of transition categories instead of a single category.

This is because a symbol e.g. 'a', '9' etc, might be valid for multiple categories. For instance 'a' is considered to be both a letter and a number in hexadecimal.

The DFSA for accepting the micro-syntax PArL is built as follows.

Let  $\mathcal{U}$  be the set of all possible characters under the system encoding (e.g. UTF-8).

The will use the following categories:

- $L := \{A, \dots, Z, a, \dots, z\}$
- $D := \{0, \dots, 9\}$
- $H := \{A, \dots, F, a, \dots, f\} \cup D$
- $S := \{\alpha \in \mathcal{U}: \alpha \text{ is whitespace}\} \setminus \{\text{LF}\}$

Note: LF refers to line-feed or as it is more commonly known '\n' i.e. new-line.

Together these categories form our alphabet  $\Sigma$ :

$$\Sigma := L \cup D \cup S \cup \{., \#, \_, (, ), [, ], \{, \}, *, /, +, -, <, >, =, !, ,, :, ;, \text{LF}\}$$

Now, the following drawing describe the transitions of the DFSA. For improved readability the DFSA has been split across multiple drawings. Hence, in each drawing initial state 0 refers to the *same* initial state (a DFSA has one and only one initial state).

Additionally, each final state is annotated with the token type it should produce.

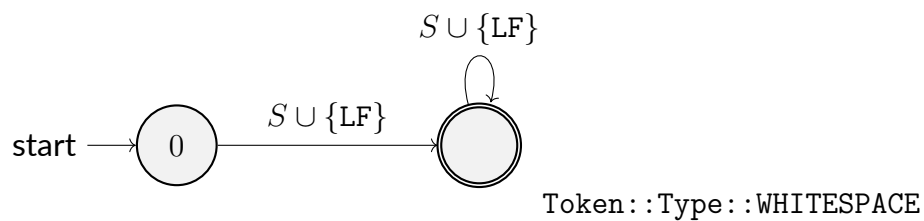


Figure 1: States & transitions for recognising whitespace

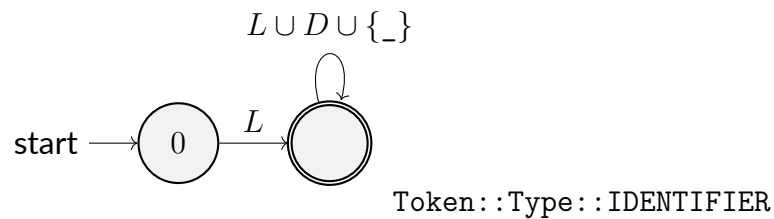


Figure 2: States & transitions for recognising identifiers/keywords

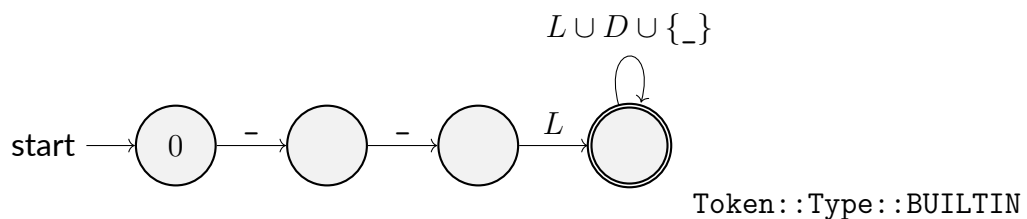


Figure 3: States & transitions for recognising builtins

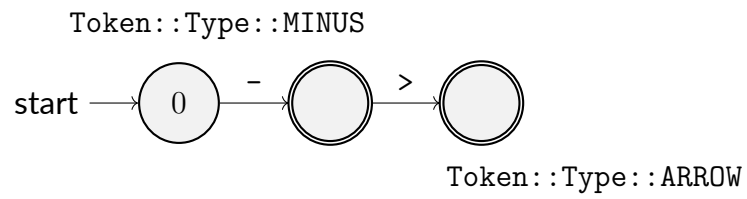


Figure 4: States &amp; transitions for recognising minus and arrow (-&gt;)

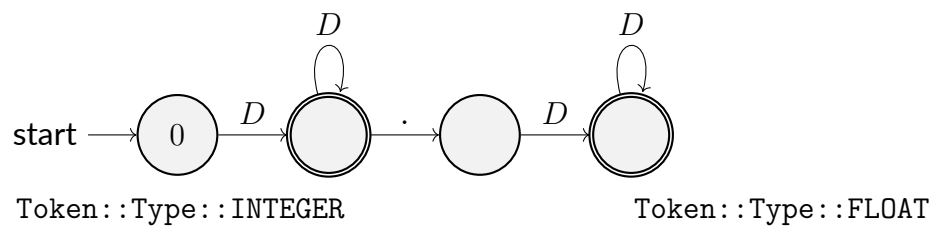


Figure 5: States &amp; transitions for recognising integers and floats

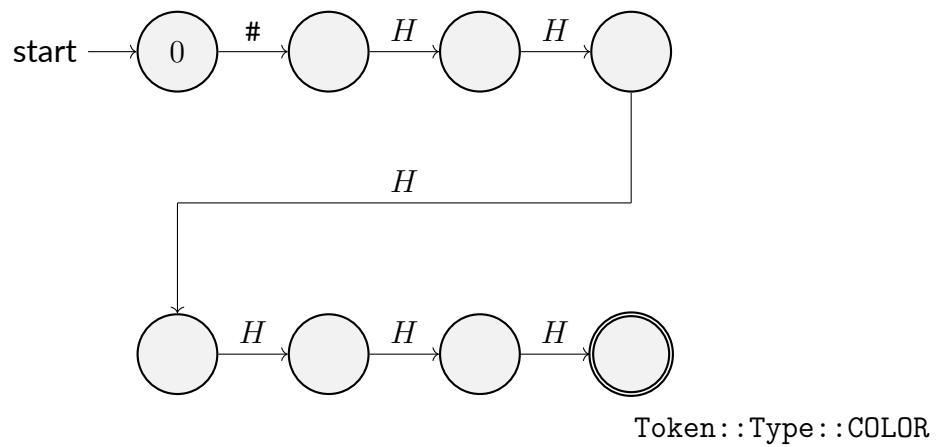


Figure 6: States &amp; transitions for recognising colours



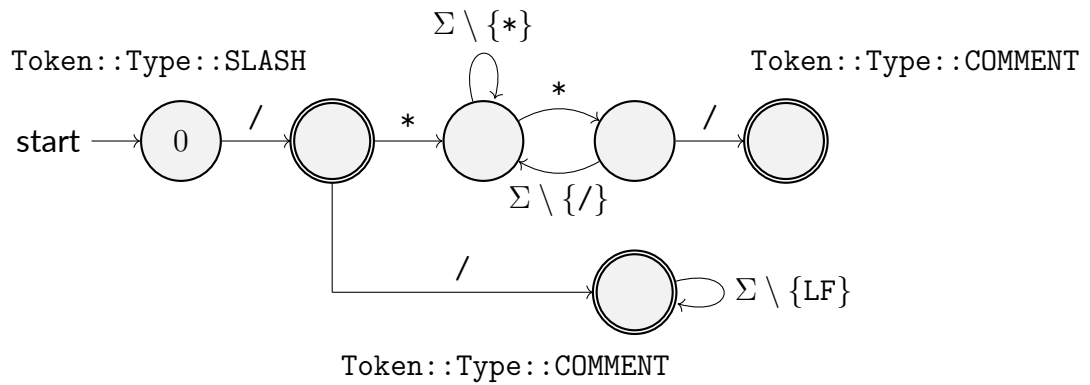


Figure 7: States &amp; transitions for recognising slashes and comments

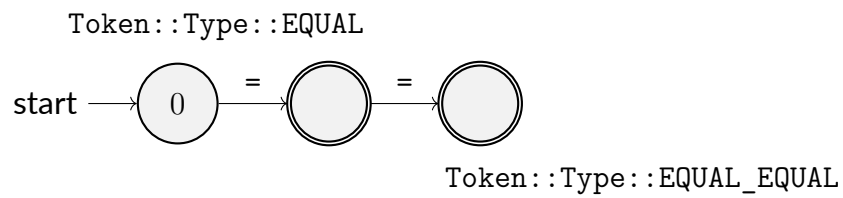


Figure 8: States &amp; transitions for assign and is equal to

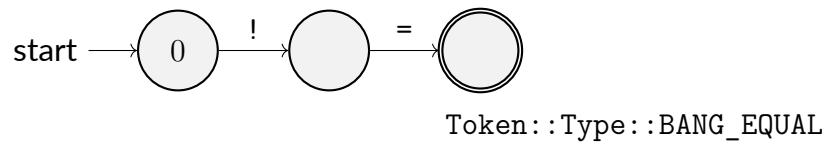


Figure 9: States &amp; transitions for not equal to

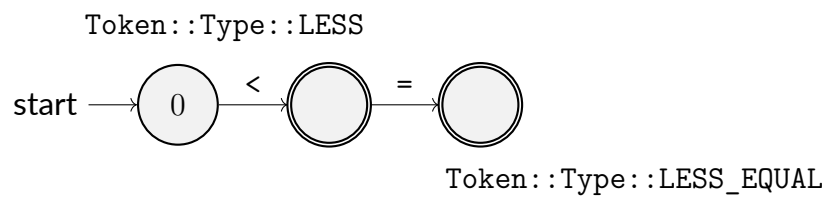


Figure 10: States &amp; transitions for less than and less than or equal to

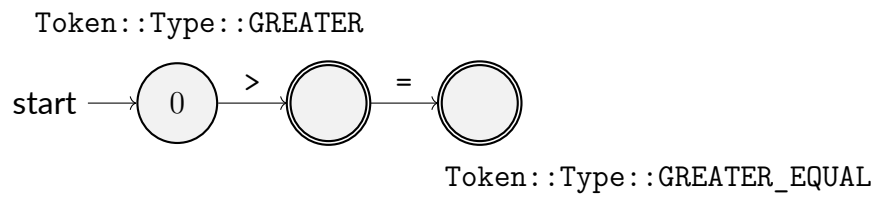


Figure 11: States &amp; transitions for greater than and greater than or equal to

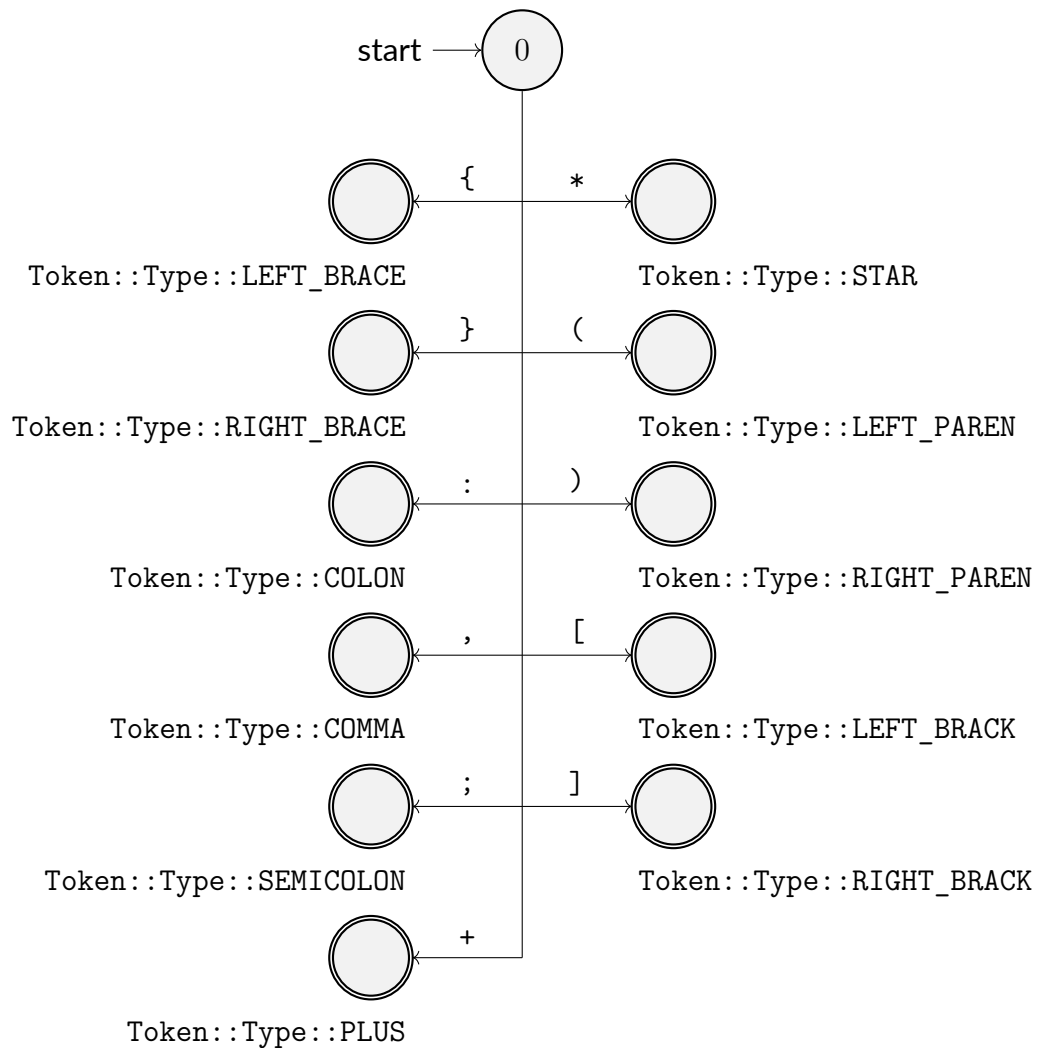


Figure 12: States &amp; transitions for single letter tokens

## The Builder & Director

Each sequence of states present is directly represented in code within the `LexerDirector` using methods provided by the `LexerBuilder`.

---

```
// "/", "//", "/* ... */"
builder.addTransition(0, SLASH, 34)
    .setStateAsFinal(34, Token::Type::SLASH)
    .addTransition(34, SLASH, 35)
    .addComplementaryTransition(35, LINEFEED, 35)
    .setStateAsFinal(35, Token::Type::COMMENT)
    .addTransition(34, STAR, 36)
    .addComplementaryTransition(36, STAR, 36)
    .addTransition(36, STAR, 37)
    .addComplementaryTransition(37, SLASH, 36)
    .addTransition(37, SLASH, 38)
    .setStateAsFinal(38, Token::Type::COMMENT);
```

---

Listing 2: Code specification of the comments in the `LexerDirector` (`lexer/LexerDirector.cpp`)

The `LexerBuilder` keeps track of these transitions using less efficient data structures such as hash maps (`std::unordered_map`) and sets (`std::unordered_set`).

Then the `build()` method processes the user defined transitions and normalises everything into a single transition table for use in a DFSA. Additionally, it also produces two other artefacts. The first is called `categoryIndexToChecker`. It is a hash map from the index of a category to a lambda function which takes a character as input and returns true or false.

The lambdas and the category indices are also registered by the user. See Listing 3 for a registration example. Additionally, the category indices although they are integers for readability they are defined as an enumeration.

---

```
.addCategory(
    HEX,
    [](char c) -> bool {
        return ('0' <= c && c <= '9') ||
               ('A' <= c && c <= 'F') ||
               ('a' <= c && c <= 'f');
    })
)
```

---

**Listing 3: Registration of the hexadecimal category checker (lexer/LexerDirector.cpp)**

The second artefact produced by the builder is also a hash map from final states to their associated token type.

The transition table is then passed onto the DFSA. And the DFSA, and the two artefacts are passed onto the Lexer class.

---

```
// create dfsa
Dfsa dfsa(
    noOfStates,
    noOfCategories,
    transitionTable,
    initialStateIndex,
    finalStateIndices
);

// create lexer
Lexer lexer(
    std::move(dfsa),
    std::move(categoryIndexToChecker),
    std::move(finalStateIndexToTokenType)
);
```

---

**Listing 4: Constructions of the Lexer (lexer/LexerBuilder.cpp)**

## The Actual Lexer

The lexer's core is as was described during the lectures and the core/main method is `simulateDFSA()`.

It also has a number of very important auxiliary methods and behavioural changes. Specifically, the `updateLocationState()`, see Listing 5, is critical for providing adequate error messages both during the current stage and for later stages. This function is called every time a lexeme is consumed allowing the lexer to keep track of where in the file it is, in terms of lines and columns.

---

```
void Lexer::updateLocationState(std::string const& lexeme) {
```

```
for (char ch : lexeme) {
    mCursor++;

    if (ch == '\n') {
        mLine++;

        mColumn = 1;
    } else {
        mColumn++;
    }
}
}
```

---

Listing 5: The `updateLocationState()` lexer method (lexer/Lexer.cpp)

Additionally, if an invalid / non-accepting state is reached the invalid lexeme is consumed and the user is warned, see Listing 6. After this the lexer, is left in a still operational state. Hence, `nextToken()` can be used again.

This is critical to provide users of the PArL compiler with a list of as many errors as possible, since it would be a bad experience to have to constantly run the PArL compiler to see the next error.

---

```
if (state == INVALID_STATE) {
    mHasError = true;

    fmt::println(
        stderr,
        "lexical error at {}:{}: unexpected "
        "lexeme '{}'",
        mLine,
        mColumn,
        lexeme
    );
} else {
    try {
        token = createToken(
            lexeme,
            mFinalStateToTokenType.at(state)
        );
    } catch (UndefinedBuiltin& error) {
        mHasError = true;
    }
}
```

```
        fmt::println(  
            stderr,  
            "lexical error at {}:{}:: {}",  
            mLine,  
            mColumn,  
            error.what()  
        );  
    }  
}
```

---

Listing 6: Error handling mechanism in the `nextToken()` lexer method (lexer/Lexer.cpp)

## Hooking up the Lexer to the Runner

The Runner class is the basic structure which connects all the stages of the compiler together together.

In this case the Runner passes in a reference to the lexer into the parser, this allows the parser to request tokens and they are computed on demand improving overall performance. Additionally, this has the benefit of allowing the parsing of larger and multiple files since, the parser is no longer limited by the amount of usable memory, since it does not need to load the whole file.

However, in this case no such optimisation is present.

---

```
Runner::Runner(bool dfsaDbg, bool lexerDbg, bool parserDbg)  
    : mDfsaDbg(dfsaDbg),  
      mLexerDbg(lexerDbg),  
      mParserDbg(parserDbg),  
      mLexer(LexerDirector::buildLexer()),  
      mParser(Parser(mLexer)) {  
}
```

---

Listing 7: The Runner constructor passes `mLexer` into the Parser constructor (runner/Runner.cpp)

## 2 | The Parser

## 2.1 | Modified EBNF

Some modifications were applied to the original EBNF. Some of the modifications were either motivated by improved user experience, a more uniform mechanism and others to reduce complexity further down the pipeline.

$\langle \text{Letter} \rangle$	::= 'A'-'Z'   'a'-'z'
$\langle \text{Digit} \rangle$	::= '0'-'9'
$\langle \text{Hex} \rangle$	::= 'A'-'F'   'a'-'f'   $\langle \text{Digit} \rangle$
$\langle \text{Identifier} \rangle$	::= $\langle \text{Letter} \rangle$ { '_'   $\langle \text{Letter} \rangle$   $\langle \text{Digit} \rangle$ }
$\langle \text{BooleanLiteral} \rangle$	::= 'true'   'false'
$\langle \text{IntegerLiteral} \rangle$	::= $\langle \text{Digit} \rangle$ { $\langle \text{Digit} \rangle$ }
$\langle \text{FloatLiteral} \rangle$	::= $\langle \text{Digit} \rangle$ { $\langle \text{Digit} \rangle$ } '.' $\langle \text{Digit} \rangle$ { $\langle \text{Digit} \rangle$ }
$\langle \text{ColorLiteral} \rangle$	::= '#' $\langle \text{Hex} \rangle$ $\langle \text{Hex} \rangle$ $\langle \text{Hex} \rangle$ $\langle \text{Hex} \rangle$ $\langle \text{Hex} \rangle$ $\langle \text{Hex} \rangle$
$\langle \text{ArrayLiteral} \rangle$	::= '[' [ $\langle \text{Epxr} \rangle$ { ',' $\langle \text{Epxr} \rangle$ } ] '
$\langle \text{PadWidth} \rangle$	::= '__width'
$\langle \text{PadHeight} \rangle$	::= '__height'
$\langle \text{PadRead} \rangle$	::= '__read' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$
$\langle \text{PadRandomInt} \rangle$	::= '__random_int' $\langle \text{Epxr} \rangle$
$\langle \text{Literal} \rangle$	::= $\langle \text{BooleanLiteral} \rangle$   $\langle \text{IntegerLiteral} \rangle$   $\langle \text{FloatLiteral} \rangle$   $\langle \text{ColorLiteral} \rangle$   $\langle \text{ArrayLiteral} \rangle$   $\langle \text{PadWidth} \rangle$   $\langle \text{PadHeight} \rangle$   $\langle \text{PadRead} \rangle$   $\langle \text{PadRandomInt} \rangle$
$\langle \text{Type} \rangle$	::= ('bool'   'int'   'float'   'color') [ '[' $\langle \text{IntegerLiteral} \rangle$ ']' ]

---

$\langle \text{SubExpr} \rangle$	::= ' (' $\langle \text{Expr} \rangle$ ')'
$\langle \text{Variable} \rangle$	::= $\langle \text{Identifier} \rangle$
$\langle \text{ArrayAccess} \rangle$	::= $\langle \text{Identifier} \rangle$ '[' $\langle \text{Expr} \rangle$ ']'
$\langle \text{FunctionCall} \rangle$	::= $\langle \text{Identifier} \rangle$ ' (' [ $\langle \text{Expr} \rangle$ {',' $\langle \text{Expr} \rangle$ }] ')'
$\langle \text{Expr} \rangle$	::= $\langle \text{LogicOr} \rangle$ ['as' $\langle \text{Type} \rangle$ ]
$\langle \text{LogicOr} \rangle$	::= $\langle \text{LogicAnd} \rangle$ {'or' $\langle \text{LogicAnd} \rangle$ }
$\langle \text{LogicAnd} \rangle$	::= $\langle \text{Equality} \rangle$ {'and' $\langle \text{Equality} \rangle$ }
$\langle \text{Equality} \rangle$	::= $\langle \text{Comparison} \rangle$ {'=='   '!=' $\langle \text{Comparison} \rangle$ }
$\langle \text{Comparison} \rangle$	::= $\langle \text{Term} \rangle$ {'<'   '<='   '>'   '>=' $\langle \text{Term} \rangle$ }
$\langle \text{Term} \rangle$	::= $\langle \text{Factor} \rangle$ {'+'   '-' $\langle \text{Factor} \rangle$ }
$\langle \text{Factor} \rangle$	::= $\langle \text{Unary} \rangle$ {'*'   '/' $\langle \text{Unary} \rangle$ }
$\langle \text{Unary} \rangle$	::= ('-'   'not') $\langle \text{Unary} \rangle$   $\langle \text{Primary} \rangle$
$\langle \text{RefExpr} \rangle$	::= $\langle \text{Variable} \rangle$   $\langle \text{ArrayAccess} \rangle$   $\langle \text{FunctionCall} \rangle$
$\langle \text{Primary} \rangle$	::= $\langle \text{Literal} \rangle$   $\langle \text{SubExpr} \rangle$   $\langle \text{RefExpr} \rangle$
$\langle \text{Program} \rangle$	::= { $\langle \text{Stmt} \rangle$ }
$\langle \text{Stmt} \rangle$	::= $\langle \text{Block} \rangle$   $\langle \text{VariableDecl} \rangle$ ';'   $\langle \text{FunctionDecl} \rangle$   $\langle \text{Assignment} \rangle$ ';'   $\langle \text{PrintStmt} \rangle$ ';'   $\langle \text{DelayStmt} \rangle$ ';'   $\langle \text{WriteBoxStmt} \rangle$ ';'   $\langle \text{WriteStmt} \rangle$ ';'   $\langle \text{ClearStmt} \rangle$ ';'   $\langle \text{IfStmt} \rangle$



	$\langle \text{ForStmt} \rangle$
	$\langle \text{WhileStmt} \rangle$
	$\langle \text{ReturnStmt} \rangle$ ';'
$\langle \text{Block} \rangle$	::= '{' { $\langle \text{Stmt} \rangle$ } '}'
$\langle \text{VariableDecl} \rangle$	::= 'let' $\langle \text{Identifier} \rangle$ ':' $\langle \text{Type} \rangle$ '=' $\langle \text{Epxr} \rangle$
$\langle \text{FormalParam} \rangle$	::= $\langle \text{Identifier} \rangle$ ':' $\langle \text{Type} \rangle$
$\langle \text{FunctionDecl} \rangle$	::= 'fun' $\langle \text{Identifier} \rangle$ '(' [ $\langle \text{FormalParam} \rangle$ {',' $\langle \text{FormalParam} \rangle$ } ] ')' '->' $\langle \text{Type} \rangle$ $\langle \text{Block} \rangle$
$\langle \text{Assignment} \rangle$	::= $\langle \text{Identifier} \rangle$ '[' $\langle \text{Epxr} \rangle$ ']' '=' $\langle \text{Epxr} \rangle$
$\langle \text{PrintStmt} \rangle$	::= '__print' $\langle \text{Epxr} \rangle$
$\langle \text{DelayStmt} \rangle$	::= '__delay' $\langle \text{Epxr} \rangle$
$\langle \text{WriteBoxStmt} \rangle$	::= '__write_box' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$
$\langle \text{WriteStmt} \rangle$	::= '__write' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$ ',' $\langle \text{Epxr} \rangle$
$\langle \text{ClearStmt} \rangle$	::= '__clear' $\langle \text{Epxr} \rangle$
$\langle \text{IfStmt} \rangle$	::= 'if' '(' $\langle \text{Expr} \rangle$ ')' $\langle \text{Block} \rangle$ ['else' $\langle \text{Block} \rangle$ ]
$\langle \text{ForStmt} \rangle$	::= 'for' '(' [ $\langle \text{VariableDecl} \rangle$ ] ';' $\langle \text{Expr} \rangle$ ';' [ $\langle \text{Assignment} \rangle$ ] ')' $\langle \text{Block} \rangle$
$\langle \text{WhileStmt} \rangle$	::= 'while' '(' $\langle \text{Expr} \rangle$ ')' $\langle \text{Block} \rangle$
$\langle \text{ReturnStmt} \rangle$	::= 'return' $\langle \text{Expr} \rangle$

### Improved Precedence

So, the minor changes which improve programmer usability are the additions of a number of other expression stages, such as  $\langle \text{LogicOr} \rangle$ ,  $\langle \text{LogicAnd} \rangle$ , etc. The main reason for the addition of such rules is to further enforce a more natural operation precedence. For example a programmer often expects that comparison operators such as  $<$  and  $>$  bind tighter than  $\text{and}$  or  $\text{or}$ , hence the compiler needs to make sure that comparison operators are executed before logical operators, and this can be enforced by the grammar itself hence the changes.

## Better Arrays

The way arrays were being implemented in the original grammar was very restrictive. Instead an approach for treating arrays as their own type and literal was taken up.

The  $\langle \text{Type} \rangle$  and  $\langle \text{Literal} \rangle$  productions were augmented to improve array support. This helped simplify the  $\langle \text{Identifier} \rangle$  (in the original EBNF),  $\langle \text{VariableDecl} \rangle$  and  $\langle \text{FormalParam} \rangle$  productions. Additionally, this opens up further support for more complicated types later on.

For example, the  $\langle \text{Type} \rangle$  can be further augmented to support more types.

```

 $\langle \text{StructField} \rangle$       ::=  $\langle \text{Identifier} \rangle$  ':'  $\langle \text{Type} \rangle$ 

 $\langle \text{Struct} \rangle$           ::= 'struct'  $\langle \text{Identifier} \rangle$  '{' { $\langle \text{StructField} \rangle$  ';' } '}'

 $\langle \text{TypeDecl} \rangle$       ::=  $\langle \text{Struct} \rangle$ 

 $\langle \text{Base} \rangle$            ::= ('bool' | 'int' | 'float' | 'color')

 $\langle \text{Array} \rangle$           ::=  $\langle \text{Type} \rangle$  '['  $\langle \text{IntegerLiteral} \rangle$  ']'

 $\langle \text{Pointer} \rangle$         ::=  $\langle \text{Type} \rangle$  '*'

 $\langle \text{Type} \rangle$            ::=  $\langle \text{Identifier} \rangle$ 
                        |  $\langle \text{Base} \rangle$ 
                        |  $\langle \text{Array} \rangle$ 
                        |  $\langle \text{Pointer} \rangle$ 

```

Note that we are indeed repeating  $\langle \text{FormalParam} \rangle$  but this is not really a problem. When it comes to specification repetition which improves clarity is “good” repetition.

Additionally, some of the ground work for these improvements has already been laid out in the internal type system see Listing 33 and Listing 9.

---

```
struct Primitive;
```

```
enum class Base {
    BOOL,
    COLOR,
    FLOAT,
```

```
    INT,  
};  
  
struct Array {  
    size_t size;  
    box<Primitive> type;  
};
```

---

Listing 8: Mechanism for internally storing types within the compiler (parl/Core.hpp)

---

```
struct Primitive {  
    template <typename T>  
    [[nodiscard]] bool is() const {  
        return std::holds_alternative<T>(data);  
    }  
  
    template <typename T>  
    [[nodiscard]] const T &as() const {  
        return std::get<T>(data);  
    }  
  
    ...  
  
    bool operator!=(Primitive const &other) {  
        return !operator==(other);  
    }  
  
    std::variant<std::monostate, Base, Array> data{};  
};
```

---

Listing 9: The Primitive structure (parl/Core.hpp)

The box type is a special type of pointer object which has value semantics that is it behaves as though it were the object it contains.

This is critical because self-referential types like the `<Array>` as described above are not easily representable within C++ since, something like Listing 10, is not allowed.

---

```
struct SelfRef;
```

```
struct Container {  
    Other other;  
    SelfRef ref;  
};  
  
struct Primitive {  
    std::variant<std::monostate, Container> data{};  
};
```

---

Listing 10: An size unbounded type in C++ (parl/Core.hpp)

This is because the C++ compiler is incapable of determining the size of said type at compile-time. Hence, a pointer for said type is required. And the pointer wrapper `box<>` allows it to be copied as though it were a value.

The `box<>` type is attributed to Jonathan Müller where he discusses the exact issue described above on his [blog](#), foonathan.

These changes to type also require additional changes to how variables are referenced in expressions hence why `<RefExpr>` was added. It also ensures that in the future more complicated referencing such as `object.something` (struct member referencing) is possible.

Finally, the `<ArrayLiteral>` has been improved to support expressions instead of just only literals.

## Removing Eye-Candy

Adding this system however, significantly increases the complexity of semantic analysis.

Because of this the slight syntax sugar of being able to have the following two conveniences `let a: int[] = [1,1,1];` and `let a: int[3] = [1];` has been dropped.

This is because such syntax not only complicates the grammar but it also significantly increases the complexity of semantic analysis.

The best way to handle this is to have a de-sugaring & type inference sub-phase before type checking, within the semantic analysis phase.

## 2.2 | Parsing & The Abstract Syntax Tree

The AST is the data structure which is produced by the parser. The nodes of the AST are in fact almost a one-to-one representation of the productions in the grammar, for example see Listing 11.

---

```
void accept(Visitor*) override;

std::unique_ptr<Expr> subExpr;
};

struct Binary : public Expr {
    explicit Binary(std::unique_ptr<Expr>, Operation, std::
unique_ptr<Expr>);
```

---

Listing 11: The FunctionCall AST node class (parl/AST.hpp)

The only significant difference in these nodes in the `Position` field. This get populated by the parser using the location of a token in the original source file. This is again critical for adequate error messaging later in the semantic analysis phase.

---

```
void accept(Visitor*) override;

const Operation op;
std::unique_ptr<Expr> expr;
};

struct Stmt : public Node {};

struct Assignment : public Stmt {
```

---

Listing 12: The Binary AST node class (parl/AST.hpp)

```
void accept(Visitor*) override;

const std::string identifier;
std::unique_ptr<Expr> index;
std::unique_ptr<Expr> expr;
```

```
};  
  
struct VariableDecl : public Stmt {
```

---

Listing 13: The Unary AST node class (parl/AST.hpp)

Apart from these the only other real difference is the use of a Binary and Unary node see Listing 12 and Listing 13, instead of a node for each type of the grammar rules discussed above in 2.1. This is because the reason for said rules was precedence and precedence within the AST is actually described by the structure of the tree itself not the node types.

Additionally, a number of the nodes override the `accept()` method specified by the pure virtual class `Node`. This is the basis for the Visitor pattern which apart from the parser is the backbone of the remaining phases.

## 2.3 | The Actual Parser

The parser is split into four main sections.

- AST Generation
- Token Buffering
- Token Matching
- Error Handling/Recovery

### Token Buffering

The parser requires access to the tokens for proper functioning. Additionally, sometimes the parser requires more than one token as the case deciding on whether an identifier is just a variable or a function call. This is quite easy to implement if the parser has available to it at initialisation all the tokens.

However, as described in 1.1. This is not the case the parser requests token on demand from lexer which it has a reference. This of course means that the machinery for handling tokens is a bit more complicated due to requiring lookahead.

To solve this issue a window-based approach was adopted. The parser has a moving buffer/window called `mTokenBuffer` and whose size is specified at compile-time using a C-style macro `#define LOOKAHEAD (2)`. The core methods for this aspect of the parser are `moveWindow()` and `nextToken()`, see Listing 14.

---

```
Token token = previous();

consume(
    Token::Type::LEFT_PAREN,
    "expected '(' after identifier"
);

std::vector<std::unique_ptr<core::Expr>> params{};

if (!peekMatch({Token::Type::RIGHT_PAREN})) {
    do {
        params.emplace_back(expr());
    } while (match({Token::Type::COMMA}));
}

consume(
    Token::Type::RIGHT_PAREN,
    "expected ')' after parameters"
);
```

---

Listing 14: The `moveWindow()` and `nextToken()` Parser methods (parser/Parser.cpp)

Additionally, note that within `nextToken()` the whitespace and comments are being explicitly ignored.

This might lead to further minor improvement. Specifically, comments can be integrated into the AST. This basically allows printing visitors or, formatting visitors properly format the code whilst still preserving any comments created by programmers.

## Token Matching

The previous methods all facilitate the more important token matching methods which are:

- `peek()`,
- `advance()`,
- `previous()`,

- `isAtEnd()`,
- `peekMatch()`,
- `match()`,
- and, `consume()`.

Arguably, the most important of these methods is `consume()`, it takes in a token type and a error message if the specified token type is not matched, see Listing 15.

---

```
template <typename... T>
void consume(
    Token::Type type,
    fmt::format_string<T...> fmt,
    T&&... args
) {
    if (check(type)) {
        advance();
    } else {
        error(fmt, args...);
    }
}
```

---

Listing 15: The `consume()` Parser methods (parser/Parser.hpp)

The main reason for the templating of such a method is to provide an easy interface for formattable strings using `fmtlib`. The main feature this library provides is the ability to specify placeholders in the string itself using `{}`. Usage of this functionality is demonstrated in the AST generator methods, see Listing 16.

---

```
Token::Type::RIGHT_PAREN,
    "expected ')' after expression"
);

std::unique_ptr<core::Block> block_ = block();
```

---

Listing 16: The Usage of the `consume()` method in the `formalParam()` generator method (parser/Parser.cpp)

The `peekMatch()` method is a simple method which returns true if at least one of the provided token types match. The main usage for this is cases where more



than different token is valid option for example such is the case for `<Comparison>` production. `match()` is a simple extension of `peekMatch()` which consumes the token if it matches.

The methods `advance()`, `previous()` and `isAtEnd()` are quite self explanatory. `advance()` moves the buffer window one step forward, `previous()` returns the last consumed token, and `isAtEnd()` checks whether or not the parser has reached an End of File token.

`peek()` is also very simple it allows the parser to see the token without consuming it. However, since the parser might need to lookahead `peek()` supports an offset. The only issue is that all calls to `peek()` have to be checked to ensure that valid access, see Listing 17. This is done using `abort_if()` function which prints an error message and aborts the program. This is very similar to a `static_assert` however, it has to be performed at runtime, see Listing 18.

---

```
}

void Parser::initWindow() {
    for (int i = 0; i < LOOKAHEAD; i++) {
        mTokenBuffer[i] = nextToken();
    }

    mPreviousToken = mTokenBuffer[0];
}
```

---

Listing 17: The `peek()` Parser method (parser/Parser.cpp)

---

```
#ifdef NDEBUG
template <typename... T>
inline void abort(fmt::format_string<T...>, T &&...) {}
}

template <typename... T>
inline void
abort_if(bool, fmt::format_string<T...>, T &&...) {}
}
#else
template <typename... T>
inline void
abort(fmt::format_string<T...> fmt, T &&...args) {
    fmt::println(stderr, fmt, args...);
}
```

```
        std::abort();
    }

    template <typename... T>
    inline void abort_if(
        bool cond,
        fmt::format_string<T...> fmt,
        T &&...args
    ) {
        if (cond) {
            abort(fmt, args...);
        }
    }
#endif
```

---

Listing 18: The abort functionality present in the codebase (parl/Core.hpp)

Note that since the usage of `abort()` and `abort_if()` is internal to the code, it only makes sense for these functions to be enabled during debug builds only. Hence, the implementation are enclosed between an `#ifdef`, `#else`, `#endif` macro. When `NDEBUG` which stand for no-debug is defined the function bodies are hollowed out allowing the C++ compiler to optimise them out.

## Error Handling/Synchronization

Error handling and synchronization is a critical part of the parser. With regards to developer productivity, having meaningful errors is extremely valuable. But apart from that being able to see all the errors in a file is also crucial. This means that a developer will waste less time re-running the compiler to find all the errors present in the source code. This process is referred to as [Synchronization](#) by the author of [Crafting Interpreters](#), Robert Nystrom. The method for synchronization in the PArL compiler was inspired by the above credited author and his usage of Exceptions as an unrolling primitive although unorthodox is very simple to implement and has in fact been used in the parser and also later stages of the compiler with success.

---

```
class SyncParser : public std::exception {};
```

...

```
template <typename... T>
void error(fmt::format_string<T...> fmt, T&&... args) {
    mHasError = true;

    Token violatingToken = peek();

    fmt::println(
        stderr,
        "parsing error at {}:{}:: {}",
        violatingToken.getPosition().row(),
        violatingToken.getPosition().col(),
        fmt::format(fmt, args...)
    );

    throw SyncParser{};
}
```

---

Listing 19: The ParseSync exception and the error() method which kickstarts the synchronization process (parser/Parser.hpp)

However, there is of course a major downside to this where synchronization happens will affect other error messages which are reported down stream. This of course has the possibility of producing false positives. However, in this case error handling is only best-effort and therefore the possibility of false positives accepted. The basic process of synchronizing is consuming as many tokens as possible until the parser reaches a token which it believes to be a good restarting point, see Listing 20.

---

```
bool Parser::check(Token::Type type) {
    if (isAtEnd()) {
        return false;
    }

    return peek().getType() == type;
}

bool Parser::peekMatch(
    std::initializer_list<Token::Type> const &types
) {
    return std::any_of(
        types.begin(),
```

```
        types.end(),
        [&](auto type) {
            return check(type);
        }
    );
}

bool Parser::match(
    std::initializer_list<Token::Type> const &types
) {
    if (peekMatch(types)) {
        advance();

        return true;
    }

    return false;
}

// NOTE: synchronization is all best effort

void Parser::synchronize() {
    while (!isAtEnd()) {
        Token peekToken = peek();

        switch (peekToken.getType()) {
            case Token::Type::SEMICOLON:
                advance();

                return;
            case Token::Type::FOR:
                /* fall through */
            case Token::Type::FUN:
```

---

Listing 20: The `synchronize()` method in the Parser class (`parser/Parser.cpp`)

Finally, the most natural choice for capturing `ParserSync` is in AST generator methods responsible for generating statement nodes, those being `program()` and `block()`.

## AST Generator Methods

Finally, the bulk of the parser is actually the generator methods which actually build the AST. These methods are not that complicated and they follow the specified grammar faithfully.

See, Listing 21 for an example of such a method and see Listing 22 for a full list of the methods.

---

```
consume(  
    Token::Type::COMMA,  
    "expected ',' after expression"  
);  
  
std::unique_ptr<core::Expr> xOffset = expr();  
  
consume(  
    Token::Type::COMMA,  
    "expected ',' after expression"  
);  
  
std::unique_ptr<core::Expr> yOffset = expr();  
  
consume(  
    Token::Type::COMMA,  
    "expected ',' after expression"  
);  
  
std::unique_ptr<core::Expr> color = expr();  
  
return make_with_pos<core::WriteBoxStmt>(  
    token.getPosition(),  
    std::move(x),  
    std::move(y),  
    std::move(xOffset),  
    std::move(yOffset),  
    std::move(color)  
);  
}  
  
std::unique_ptr<core::IfStmt> Parser::ifStmt() {  
    consume(  
        Token::Type::IF,
```

---

"expected 'if' at start of if statement"

---

Listing 21: The ifStmt() node generator method in the Parser class (parser/Parser.cpp)

---

```
std::unique_ptr<core::Type> type();

std::unique_ptr<core::Program> program();
std::unique_ptr<core::Stmt> statement();
std::unique_ptr<core::Block> block();
std::unique_ptr<core::VariableDecl> variableDecl();
std::unique_ptr<core::Assignment> assignment();
std::unique_ptr<core::PrintStmt> printStatement();
std::unique_ptr<core::DelayStmt> delayStatement();
std::unique_ptr<core::WriteBoxStmt> writeBoxStatement();
std::unique_ptr<core::WriteStmt> writeStatement();
std::unique_ptr<core::ClearStmt> clearStatement();
std::unique_ptr<core::IfStmt> ifStmt();
std::unique_ptr<core::ForStmt> forStmt();
std::unique_ptr<core::WhileStmt> whileStmt();
std::unique_ptr<core::ReturnStmt> returnStmt();
std::unique_ptr<core::FunctionDecl> functionDecl();
std::unique_ptr<core::FormalParam> formalParam();

std::unique_ptr<core::PadWidth> padWidth();
std::unique_ptr<core::PadHeight> padHeight();
std::unique_ptr<core::PadRead> padRead();
std::unique_ptr<core::PadRandomInt> padRandomInt();
std::unique_ptr<core::BooleanLiteral> booleanLiteral();
std::unique_ptr<core::ColorLiteral> colorLiteral();
std::unique_ptr<core::FloatLiteral> floatLiteral();
std::unique_ptr<core::IntegerLiteral> integerLiteral();
std::unique_ptr<core::ArrayLiteral> arrayLiteral();
std::unique_ptr<core::SubExpr> subExpr();
std::unique_ptr<core::Variable> variable();
std::unique_ptr<core::ArrayAccess> arrayAccess();
std::unique_ptr<core::FunctionCall> functionCall();

std::unique_ptr<core::Expr> expr();
std::unique_ptr<core::Expr> logicOr();
std::unique_ptr<core::Expr> logicAnd();
std::unique_ptr<core::Expr> equality();
```

```
std::unique_ptr<core::Expr> comparison();  
std::unique_ptr<core::Expr> term();  
std::unique_ptr<core::Expr> factor();  
std::unique_ptr<core::Expr> unary();  
std::unique_ptr<core::Expr> primary();
```

---

Listing 22: The main body of methods in the Parser class (parser/Parser.hpp)

## 2.4 | Pretty? Printing

### Using Parser in the Runner

The parser is used in the runner and assuming that the parser does not encounter any errors and the `mParserDbg` flag is set it can be used to print the AST using a `PrinterVisitor`.

---

```
mParser.parse(source);  
  
if (mLexer.hasError() || mParser.hasError()) {  
    return;  
}  
  
std::unique_ptr<core::Program> ast = mParser.getAst();  
  
if (mParserDbg) {  
    debugParsing(ast.get());  
}
```

---

Listing 23: The parse segment of the `run()` method in the Runner class (runner/Runner.cpp)

Calling the produced PArL binary with the `-p` flag will set the `mParserDbg`, see Figure 13 and Figure 14.

```
> cat testing/test9.parl
fun AverageOfTwo(x: int, y: int) -> float {
  let t0: int = x + y;
  if (t0 > 10) {
    return 10;
  }
  let t1: float = t0 / 2 as float;
  return t1;
}
```

Figure 13: cat of the test9.parl



```

> ./run.sh -p testing/test9.parl
Parser Debug Print
Program =>
  Func Decl AverageOfTwo =>
    Formal Param x =>
      int
    Formal Param y =>
      int
    float
    {
      let t0 :
        int
        Binary Operation + =>
          Variable x
          Variable y
      If =>
        Binary Operation > =>
          Variable t0
          int 10
        {
          Return =>
            int 10
        }
      let t1 :
        float
        Binary Operation / =>
          Variable t0
          int 2
        as
          float
      Return =>
        Variable t1
    }
semantic error at 4:9:: incorrect return type in function AverageOfTwo

```

Figure 14: The AST generated by the syntactically correct program in test9.parl

```

void Runner::debugParsing(core::Program* program) {
    fmt::println("Parser Debug Print");

    PrinterVisitor printer;

    program->accept(&printer);
}

```

---

Listing 24: The `debugParsing()` method in the `Runner` class (`runner/Runner.cpp`)

The `PrinterVisitor` used in Listing 24 is a specialization of the pure virtual `Visitor` class in `parl/Visitor.hpp`.

---

```
...

struct ClearStmt;
struct Block;
struct FormalParam;
struct FunctionDecl;
struct IfStmt;
struct ForStmt;
struct WhileStmt;
struct ReturnStmt;
struct Program;

class Visitor {
public:
    virtual void visit(Type*) = 0;
    virtual void visit(Expr*) = 0;
    virtual void visit(PadWidth*) = 0;
    virtual void visit(PadHeight*) = 0;
    virtual void visit(PadRead*) = 0;
    virtual void visit(PadRandomInt*) = 0;
    virtual void visit(BooleanLiteral*) = 0;
    virtual void visit(IntegerLiteral*) = 0;
    virtual void visit(FloatLiteral*) = 0;
};

...
```

---

Listing 25: A segment of the pure virtual `Visitor` class (`parl/Visitor.hpp`)

Due to the way C++ handles symbols, the classes which the visitor can visit must be forward declared manually, see Listing 25. If no such forward declaration is made, the compiler will complain about the `Visitor` and the AST classes being cyclically dependent.

Additionally, any inheriting visitor such as the `PrinterVisitor` can hold state. In fact, this is where the true power of visitors arises. Being able to hold state

means that complex computations can be carried out on the AST. For example the `PrinterVisitor` although simple makes use of a single variable `mTabCount`, see Listing 26, which it uses to affect how much indentation should be used in printing, allowing us to visualise the AST.

---

```
void PrinterVisitor::visit(core::PadRead *expr) {
    print_with_tabs("__read =>");
    mTabCount++;
    expr->x->accept(this);
    expr->y->accept(this);
    mTabCount--;

    expr->core::Expr::accept(this);
}
```

---

Listing 26: The `visit(core::PadRead*)` method in the `PrinterVisitor` (`parser/PrinterVisitor.cpp`)

## 3 | Semantic Analysis

### 3.1 | Phases & Design

As described in 2.1 the Semantic Analysis phase should be split into a number of sub-phases specifically: symbol resolution, de-sugaring/type inference and type checking.

However, these steps can be combined together into a single step phase. There are benefits and downsides to both approaches. The main benefit of using the first approach is a reduction in algorithmic complexity. The logic can be separated into different phases with ease and information from one phase can be propagated to another. The downside of this approach is that it actually adds more code for maintenance, since each individual sub-phase will probably need to be implemented as its own visitor. The other downside is error management. If an error occurs in a particular phase since said phase is completely disjoint from the phases after it a mechanism for propagation needs to be devised which again further increases complexity. Of course by the very nature of this argument the monolithic approach does not suffer for the issues that the sub-phases approach has. However, it significantly increases complexity since all sub-phases are being done in a single phase. The other more glaring issue is the fact that it is much more difficult to resolve symbols before-hand.

You would want to do so to allow for the location of function declarations in code to not effect resolution, that is, a function can be called before it is referenced.

Unfortunately, due to the fact that compiler development was an organic process and not too much time was spend on deliberation the semantic analysis phase became monolithic, and hence it suffers from the issues described above. Of course, this means location agnostic function declaration are currently *not* supported by the compiler.

## 3.2 | The Environment Tree

Some terminology is required to properly describe the number of structures which will be used in this section. The following terms are critical and need to be differentiated properly:

- SymbolTable
- SymbolTableStack
- Environment
- EnvStack
- RefStack

During semantic analysis there is a need for specific data structures which facilitate scoping rules, type checking, etc.

The most basic data structure which achieves this, is a single SymbolTable. A symbol table in its simplest form is a wrapper around a hash map, whose keys are identifiers and values are Symbols or any structure which is capable of storing variable and function signatures.

This approach is quite limiting since it restricts developers and user of the language to a single global scope. Therefore, this structure is not a sufficient solution for most toy-languages let alone production ready languages.

A better solution is using a stack of symbol tables. This allows symbol tables to shadow each other allowing for the reuse of identifiers. Additionally, this further opens up the possibility of implementing modules at the language-level. The likelihood that names will be reused across different modules is quite high and being able to scope modules so they do not interfere with global scope and each other is necessary for any sufficiently large project.

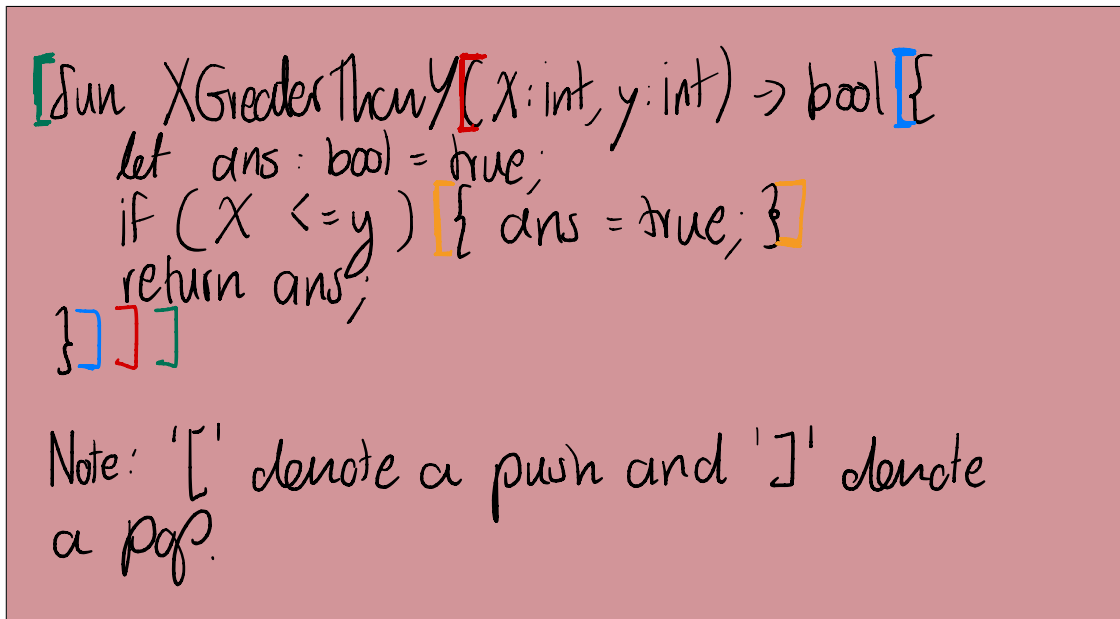


Figure 15: A PArL function with annotated scopes

The basic premise of using a symbol table stack is described in Algorithm 1 and Figure 16. A new symbol table, also referred to as a scope, is pushed onto a stack only for specific types of AST nodes. The node is then processed, where “processing” often times means considering all the sub-children of the node, and when processing terminates the scope is popped. In the context of a purely stack based implementation this implies that, scopes are only temporary that is they are lost after being popped.

---

**Algorithm 1:** Basic description of SymbolTableStack usage
 

---

**Data:**  $N$  the AST node,  $S$  the symbol table stack

**begin**

**if**  $N$  opens a scope **then**

    PushScope( $S$ );

**end**

  ProcessNode( $N$ );

**if**  $N$  opens a scope **then**

    PopScope( $S$ );

**end**

**end**

---

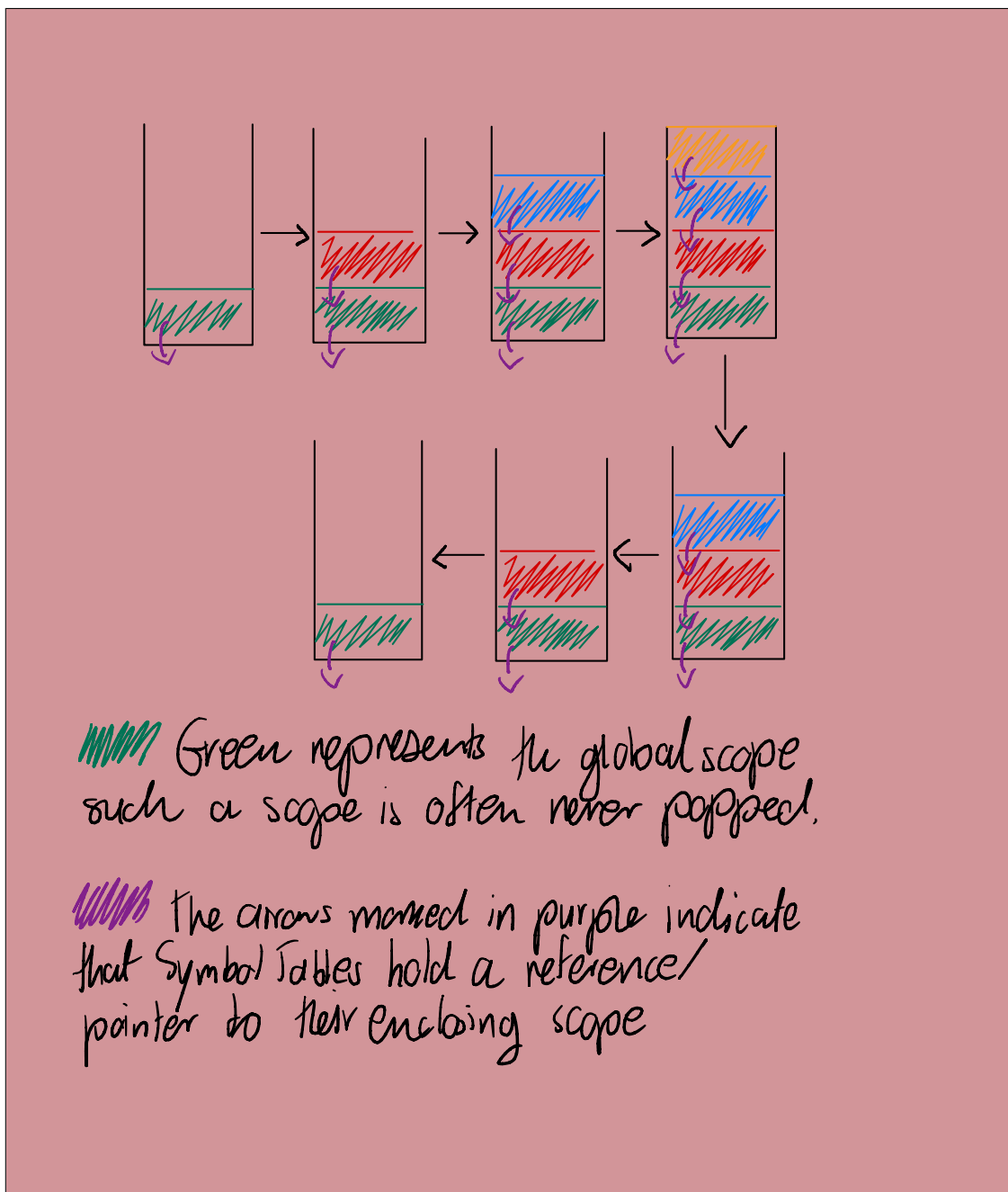


Figure 16: Behaviour of the SymbolTableStack when considering the code in Figure 15

The main advantage of using this approach is that the worst case memory usage

is  $O(DI)$  where  $D$  is the length of a longest chain of open scopes and  $I$  is the size of largest number of variable declarations within a scope.

However, due to the temporary nature of this approach it does not lend itself well to a multi-phase approach. This is because at each phase SymbolTables have to be regenerated increase the complexity of each individual phase.

Instead a different approach shall be used. This approach uses an [Environment Tree](#) and it has also been inspired by [Crafting Interpreters](#).

The notion of an Environment is essentially, the same as a SymbolTable. The only significant change is the addition of vector of unique pointers to other environments, see Listing 27.

---

```
class Environment {
public:
    enum class Type {
        GLOBAL,
        IF,
        ELSE,
        FOR,
        WHILE,
        FUNCTION,
        BLOCK
    };

    void addSymbol(
        std::string const& identifier,
        Symbol const& Symbol
    );
    [[nodiscard]] std::optional<Symbol> findSymbol(
        std::string const& identifier
    ) const;
    Symbol& getSymbolAsRef(std::string const& identifier);

    [[nodiscard]] Environment* getEnclosing() const;
    void setEnclosing(Environment* enclosing);

    [[nodiscard]] Type getType() const;
    void setType(Type type);

    [[nodiscard]] std::optional<std::string> getName(
    ) const;
```

```
void setName(const std::string& name);

std::vector<std::unique_ptr<Environment>>& children();

[[nodiscard]] bool isGlobal() const;

[[nodiscard]] size_t getIdx() const;
void incIdx();
void incIdx(size_t inc);

[[nodiscard]] size_t getSize() const;
void setSize(size_t size);

private:
    std::unordered_map<std::string, Symbol> mMap{};
    Type mType{Type::GLOBAL};
    std::optional<std::string> mName{};
    Environment* mEnclosing{nullptr};
    std::vector<std::unique_ptr<Environment>> mChildren{};
    size_t mSize{0};
    size_t mIdx{0};
};
```

---

Listing 27: The Environment class with mChildren highlighted (backend/Environment.hpp)

The additional fields present in the Environment class are there to cater for the needs of each of the phases.

Specifically, mType and mName are used during semantic analysis and mSize and mIdx are used during code generation.

Additionally, it is preferable if the interface with which the phases interact with the Environment tree is identical to that of a SymbolTableStack, that is the same push() and pop() methods are used.

Now to facilitate this another two classes need to be defined EnvStack and RefStack. The main difference between EnvStack and RefStack is that the purpose of EnvStack is to construct the Environment tree whilst RefStack traverses the environment tree. Due to this RefStack is only required during the symbol resolution. However, in current implementation since symbol resolution is combined with type checking, the Environment tree is generated at the same time as it is being type checked.



Additionally, in both an `EnvStack` and a `RefStack`, the creation and traversal of the Environment tree are managed via the `push()` and `pop()` methods.

In an `EnvStack` `push()` works by creating a new environment, setting the enclosing environment to the current environment and changing the current environment to to the new environment, see Listing ?? . `pop()` makes use of the `mEnclosing` field in the current environment and just sets the current environment to the enclosing environment essentially moving up an environment, see Listing ??.

---

```
EnvStack& EnvStack::pushEnv() {
    auto& ref = mCurrent->children().emplace_back(
        std::make_unique<Environment>()
    );

    ref->setEnclosing(mCurrent);

    mCurrent = ref.get();

    return *this;
}
```

---

Listing 28: The `pushEnv()` method in the `EnvStack` class (analysis/EnvStack.cpp)

---

```
EnvStack& EnvStack::popEnv() {
    mCurrent = mCurrent->getEnclosing();

    return *this;
}
```

---

Listing 29: The `popEnv()` method in the `EnvStack` class (analysis/EnvStack.cpp)

The implementations in `RefStack` are however a bit more involved. This is because the program has to be able to perform a step-wise left-first depth-first traversal. The best way to do this is to make use of a stack and the algorithms in Listing 30 and Listing 31.

---

```
RefStack& RefStack::pushEnv() {
    size_t index = mStack.top();

    mStack.top()++;
}
```

---

```
mStack.push(0);

mCurrent = mCurrent->children()[index].get();

return *this;
}
```

---

Listing 30: The `pushEnv()` method in the `RefStack` class (`ir_gen/RefStack.cpp`)

---

```
RefStack& RefStack::popEnv() {
    mStack.pop();

    mCurrent = mCurrent->getEnclosing();

    return *this;
}
```

---

Listing 31: The `popEnv()` method in the `RefStack` class (`ir_gen/RefStack.cpp`)

See Figure 17, for the initial segments of a traversal. The important thing to note here is that the correctness of the traversal entirely depends on the usage of `push()` and `pop()`. If used incorrectly the traversal might be incorrect or the program might crash. Of course, this is only an implementation detail that is if used correctly by the compiler developer no such issue should occur.

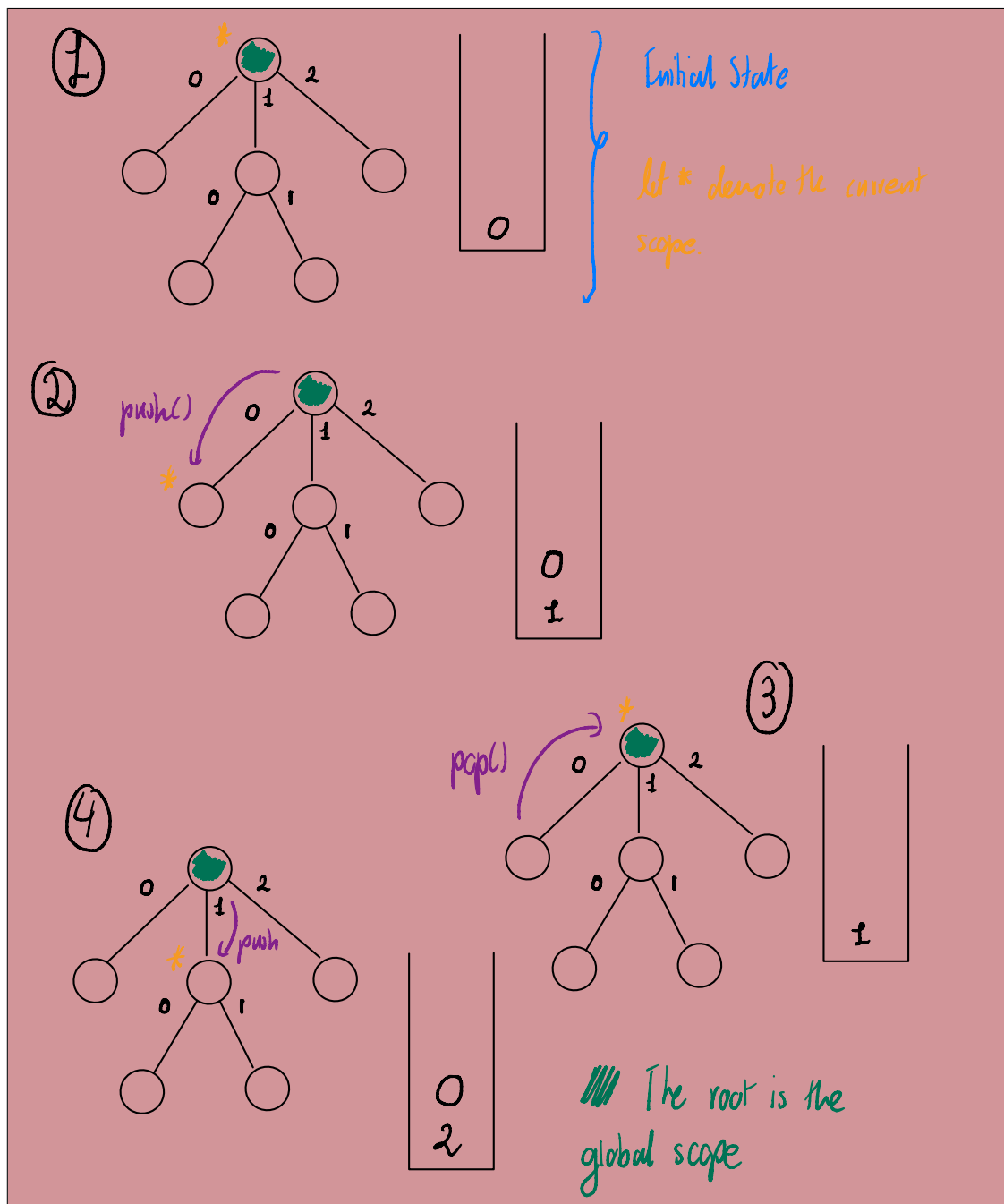


Figure 17: The initial steps of a traversal performed by RefStack

In the current implementation of PArL three types of nodes control scopes, these are: Blocks, ForStmts and FunctionDecls. There is also a slight particularity in

that, `ForStmts` and `FunctionDecls` as per the grammar actually, cause two scopes to open. This really depends on whether shadowing of, generally the iterator in a for-loop and the function parameters in the language should be allowed or not. For simplicities sake in PArL a second scope is being opened. However there is no particular consensus, for example in C++ shadowing of function parameters or declarations in for-loops is not allowed since they create a single scope, but Rust allows shadowing for the same variable in the same scope without any problems.

### 3.3 | Semantic Analysis

Having the necessary data structures in place the discussion will now revolve around Semantic Analysis i.e. the process of making sure the meaning of the program is correct.

There are four main areas of interest which shall be discussed: the symbol type, symbol registration/search, type checking and a unique subset of type checking, return value type verification.

### 3.4 | The Symbol Type

Within PArL functions are special they are not treated as values i.e. they cannot be assigned to variables, they cannot be moved around etc. This means that the `Symbol` type is actually larger than the `Primitive` type referenced in 2.1. It has to support both variables and functions, hence the `Symbol` is truly the union of two different types, `VariableSymbol` and `FunctionSymbol`.

---

```
struct VariableSymbol {
    size_t idx{0};
    core::Primitive type;

    ...
};

struct FunctionSymbol {
    size_t labelPos{0};
    size_t arity{0};
    std::vector<core::Primitive> paramTypes;
    core::Primitive returnType;

    ...
};
```

---

Listing 32: Definitions of `VariableSymbol` and `FunctionSymbol` (backend/Symbol.hpp)

Again the additional fields present, see Listing 32, within said symbols helps reduce complexity further down the pipeline especially in code generation.

Additionally, it is important to note that only primitives are comparable. Of course being able to compare the full symbols or adding said logic in the `Symbol` class goes against separation of behaviour and state.

### Symbol Registration & Search

A `Symbol` in PArL is an attribute associated with an identifier (a string). One of the main jobs of semantic analysis is to solidify these associations from the AST into the current Environment. In PArL the only nodes which are meant to cause a registration are the `Variable Declaration` node and the `Function Declaration` node.

The only important considerations for registration are that a identifier is never registered twice with in the same scope and that function identifiers are globally unique, that is they cannot be shadowed.

---

```
Environment *enclosingEnv = env->getEnclosing();

for (;;) {
    if (enclosingEnv == nullptr)
        break;

    std::optional<Symbol> identifierSignature =
        enclosingEnv->findSymbol(param->identifier);

    if (identifierSignature.has_value() &&
        identifierSignature->is<FunctionSymbol>()) {
        error(
            param->position,
            "redeclaration of {}(...) as a "
            "parameter",
            param->identifier
        );
    }
}
```

```
        enclosingEnv = enclosingEnv->getEnclosing();
    }

    env->addSymbol(param->identifier, {type});
}

void AnalysisVisitor::visit(core::FunctionDecl *stmt) {
    bool isGlobalEnv = mEnvStack.isCurrentEnvGlobal();

    if (!isGlobalEnv) {
        error(
            stmt->position,
            "function declaration {}(...) is not "
            "allowed "
            "here",
            stmt->identifier
        );
    }
}
```

---

Listing 33: The `visit(FormalParam *)` method in the `AnalysisVisitor` class (`analysis/AnalysisVisitor.cpp`)

Searching happens when a identifier in a none-declaration node is met. The procedure for searching basically requires querying the current environment for the symbol associated with the found identifier. If nothing is found climb to the enclosing state and repeat the process. If the root node or *terminating node* is reached and no associated symbol is found that means that the identifier is undefined.

Note that the specification of a *terminating node* (a `Environment*`) is important. This is because it allows the analyser to restrict access to outer scopes when necessary. The main example of this is function declarations. In PArL functions are (for the most) part pure, that is they do not produce side-effect at a language level. To enforce this, functions are not allowed to access any state outside of their own scope. The only exception to this fact is the usage of built-ins like `__write` which directly effect the simulator.

This is where one of the extra fields mentioned in 3.2, specifically `mType` is used because it allows scopes to be of different types and this is necessary to find function scopes. For example in Listing 34 `findEnclosingEnv()` is being used to find the closest enclosing function scope and if no such scope is found i.e. the optional is empty. Then searching either terminates on the stopping scope which

is either global or the enclosing function scope.

---

```
void AnalysisVisitor::visit(core::Variable *expr) {
    std::optional<Environment *> stoppingEnv =
        findEnclosingEnv(Environment::Type::FUNCTION);

    std::optional<Symbol> symbol{findSymbol(
        expr->identifier,
        stoppingEnv.has_value() ? *stoppingEnv
                               : mEnvStack.getGlobal()
    )};

    if (!symbol.has_value()) {
        error(
            expr->position,
            "{} is undefined",
            expr->identifier
        );
    }

    ...
}
```

---

Listing 34: The visit(Variable \*) method in the AnalysisVisitor class (analysis/AnalysisVisitor.cpp)

## 3.5 | Type Checking

This is the more tedious aspect of semantic. Each individual case where the types affect the behaviour of program or are simply not allowed have to be considered. For this the AnalysisVisitor has the field mReturn which is essentially used as the return of a visit, see Listing 35.

---

```
void AnalysisVisitor::visit(core::Binary *expr) {
    expr->left->accept(this);
    auto leftType{mReturn};

    expr->right->accept(this);
    auto rightType{mReturn};

    ...
}
```

---

Listing 35: A part of the `visit(Binary *)` method in the `AnalysisVisitor` class (`analysis/AnalysisVisitor.cpp`)

In terms of design choice a strict approach was chosen, that is no implicit casting takes place. Operations are not capable of operating with two distinct types, for example `5 / 2.0` is not allowed. Instead operations which should support multiple types have been overloaded (later on in code generation). For example, `__print 5 / 2` outputs `2`, `__print 5 / 2.0` is a type error and `__print 5.0 / 2.0` outputs `2.5`.

---

```
...

case core::Operation::ADD:
case core::Operation::SUB:
    if (leftType.is<core::Array>()) {
        error(
            expr->position,
            "operator {} is not defined on array "
            "types",
            core::operationToString(expr->op)
        );
    }

    if (leftType != rightType) {
        error(
            expr->position,
            "operator {} expects both operands to "
            "be of same type",
            core::operationToString(expr->op)
        );
    }

    mReturn = leftType;
    break;

...
```

---

Listing 36: Another part of the `visit(Binary *)` method in the `AnalysisVisitor` class (`analysis/AnalysisVisitor.cpp`)



Of course the remainder of the `AnalysisVisitor` is essentially, type checking in the form of described in Listing 36.

### Return Validation

---

```
...

std::optional<Environment *> optEnv =
    findEnclosingEnv(Environment::Type::FUNCTION);

if (!optEnv.has_value()) {
    error(
        stmt->position,
        "return statement must be within a "
        "function block"
    );
}

...
```

---

Listing 37: Checking whether return statement is inside a function declaration in the `visit(ReturnStmt *)` method in the `AnalysisVisitor` class (`analysis/AnalysisVisitor.cpp`)

The only non-trivial check is that of a return statement. First of all return statements must be enclosed within functions, otherwise it is a semantic error, see Listing 37.

---

```
...

Environment *env = *optEnv;

std::string enclosingFunc = env->getName().value();

auto funcSymbol = env->getEnclosing()
    ->findSymbol(enclosingFunc)
    ->as<FunctionSymbol>();

if (exprType != funcSymbol.returnType) {
    error(
        stmt->position,
```

```

        "incorrect return type in function {}",
        enclosingFunc
    );
}

...

```

Listing 38: Checking whether the return expression has the same type as the function return type in the `visit(ReturnStmt *)` method in the `AnalysisVisitor` class (`analysis/AnalysisVisitor.cpp`)

Additionally, they must have the same return type as the enclosing function, see Listing 38, and finally and most importantly, **a function must return in all possible branches**. Additionally, the implementation **should be capable of recognising unreachable areas of code**.

- ! Not stopping when an unreachable code segment is met could result in an invalid assessment of whether a function returns from all branches or not.

During the first iteration of the compiler, the mechanism for checking the return type was embedded directly into the semantic analysis. However, there is a slight problem with this approach. If an unreachable segment of code is met to avoid the issue mentioned above the remainder of the unreachable code will not be type checked. Of course this is a problem since such behaviour means the compiler can ignore semantically incorrect code and still produce VM instructions.

To solve this issue a separate visitor `ReturnVisitor` was created to handle checking that functions return from all possible branches. The separation facilitates checking after type checking is complete circumventing any of the discussed issues.

The mechanism which `ReturnVisitor` uses to establish whether a function returns or not is using a class member `mBranchReturns`.

This member is a boolean and it is explicitly set to true only by a return statement. Function declarations reset the `mBranchReturns` to false and other statements which do not exhibit branching do not modify the value of `mBranchReturns`. This means that the only statements which effect `mBranchReturns` are if-else statements, for-loops and while-loops, since they branch.

For-loops and while-loops reset `mBranchReturns` to false after their block of statements. This is because entry into the body of the loop (both while and for) is conditional and hence a program might never enter into the body of the loop. Hence,

it is not enough to return from the body of the loop.

---

```
void ReturnVisitor::visit(core::IfStmt *stmt) {
    stmt->thenBlock->accept(this);
    bool thenBranch = mBranchReturns;

    bool elseBranch = false;

    if (stmt->elseBlock) {
        stmt->elseBlock->accept(this);
        elseBranch = mBranchReturns;
    }

    mBranchReturns = thenBranch && elseBranch;
}
```

---

Listing 39: The `visit(IfStmt *)` method in the `ReturnVisitor` class (`analysis/ReturnVisitor.cpp`)

If-else statements are different in that if the ‘then’ block returns and the ‘else’ block returns then the whole if-else statement returns otherwise the if-else statement does not return. This is essentially a logical and between the results of the visitor for the ‘then’ and ‘else’ blocks, see Listing 39.

---

```
void ReturnVisitor::visit(core::Block *block) {
    size_t i = 0;

    for (; i < block->stmts.size(); i++) {
        block->stmts[i]->accept(this);

        if (mBranchReturns) {
            i++;

            break;
        }
    }

    if (i < block->stmts.size()) {
        core::Position start = block->stmts[i]->position;

        warning(
            start,
```

```
        "statements starting from line {} upto line {} "  
        "are unreachable",  
        start.row(),  
        block->position.row()  
    );  
}  
  
for (; i < block->stmts.size(); i++) { block->stmts.pop_back(); }  
}
```

---

Listing 40: The `visit(Block *)` method in the `ReturnVisitor` class (`analysis/ReturnVisitor.cpp`)

Finally, in the block statements if one of the statements inside manages to successfully return, the remaining statements will never be reached and hence the user can be notified and the statements can be dropped, see Listing 40.

The main

To do this

This minor change removes the need to throw away Environments and they can be passed on from one phase to another.

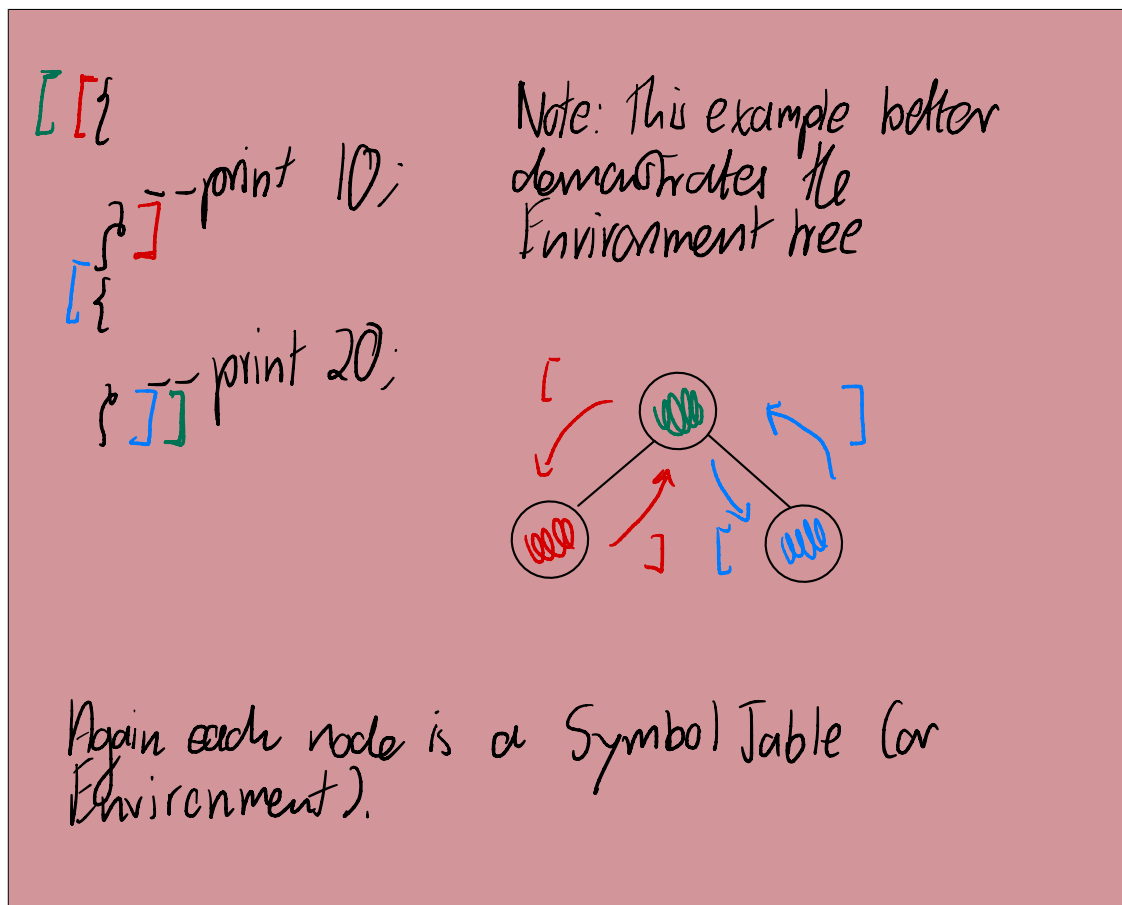


Figure 18: Construction of an Environment tree

The Environment tree is a data structure

- Describe Symbol Table.
- Describe each element of the symbol table.
- Describe the tree approach.
- attribute the tree approach to Robert Nystrom as well.
- describe the main advantage/disadvantage of using a environment approach.
- describe the main advantage/disadvantage of using a symbol stack approach.

## 4 | Code Generation

Having completed semantic analysis, the compiler can proceed to code generation.

### 4.1 | Passing on the Environment Tree

Before, proceeding onto to the details of code generation it is important to note that here is where all the upfront work that was put into making Environment trees, the data structure of choice for managing scopes, pays off. This is because there is no need to re-establish the types of any of the variables, at any scope in a program.

---

```
std::unique_ptr<Environment> environment =
    mAnalyser.getEnvironment();

ReorderVisitor reorder{};

reorder.reorderAst(ast.get());

if (mParserDbg) {
    debugParsing(ast.get());
}

reorder.reorderEnvironment(environment.get());

GenVisitor gen{environment.get()};
```

---

Listing 41: Getting the global environment from the AnalysisVisitor and passing it on for reordering and code generation (runner/Runner.cpp)

### The TypeVisitor Class

Having said that there is still the need to recompute the types of expressions since the type of a composite expression e.g.  $2 + 3$  cannot be stored in an Environment. To handle this another visitor, this time called TypeVisitor was created.

**TODO** In the event that type checking/generation is completely delegated to the TypeVisitor future version of compiler can drop usage of Environment trees in favour of re-computing types as requested. Additionally, such a change improves memory usage in exchange for time as described in 3.2.

---

```

void TypeVisitor::visit(core::FunctionCall *expr) {
    std::optional<Symbol> symbol{
        findSymbol(expr->identifier, mRefStack.getGlobal())
    };

    mReturn = symbol->as<FunctionSymbol>().returnType;

    expr->core::Expr::accept(this);
}

```

---

Listing 42: The `visit(FunctionCall *)` method in the `TypeVisitor` class (`ir_gen/TypeVisitor.cpp`)

Currently, a `TypeVisitor` is given a `RefStack` for environment access and its main capability is re-calculating the type of **already** type-checked expressions, see Listing 42.

! Usage of a `TypeVisitor` on a non-type-checked expression will at best crash and at worst, still work, leading to undefined behaviour.

## 4.2 | Reordering

Due to the fact that the PArDis VM Simulator expects all functions to be defined before the `.main` segment, which in the case of PArL is everything not in a function declaration, a solution needs to be devised to ensure that this order requirement is satisfied.

A solution to this problem can be achieved by reordering the AST and Environments. This is because function declarations can only exist in global scope, and hence moving them such that they are the first statements, which are encountered in a program is a sufficient constraint to ensure code generation as required by the VM.

This approach required two visitors, named `IsFunctionVisitor` and `ReorderVisitor`, and a function `reorderEnvironment()` (the choice to embed it into the `ReorderVisitor` was arbitrary).

The `IsFunctionVisitor` is quite self-explanatory it is a very simple visitor with an exposed method `check()` which returns true only if a node is a function declaration.

TODO

Currently the compiler has support for named scopes (see backend/Environment.cpp, lines 25 and 34). With a specially designed naming scheme to uniquely identify each scope, support for nested functions can be enabled by mangling function names with scope names and lifting function declarations into global scope.

The basic approach for solving this problem

## 4.3 | Function Declarations

Properly Closing Scopes

## 4.4 | Scopes & Variable Declarations

## 4.5 | Expressions

Type Specific VM Code

Accessing Variables

Function Calls

## 4.6 | Branching

# 5 | Attributions

- Sandro Spina for the brilliant description of table-driven lexers
- Robert Nystrom and his great book Crafting Interpreters for a great outline for parsing and error recovery/management for languages which support exceptions