
Compiler Theory and Practice

Coursework

Juan Scerri
123456A

April 1, 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 AST & Parsing	10
2.1 Modified EBNF	11

§ 3	Attributions	14
-----	------------------------	----

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 <code>updateLocationState()</code> lexer method (lexer/Lexer.cpp)	8
6	Error handling mechanism in the <code>nextToken()</code> lexer method (lexer/Lexer.cpp)	9
7	The Runner constructor passes <code>mLexer</code> into the Parser constructor (runner/Runner.cpp)	10



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.

```
14 class Dfsa {
15     public:
16         Dfsa(
17             size_t noOfStates,
18             size_t noOfCategories,
19             std::vector<std::vector<int>> const&
20                 transitionTable,
21             int initialState,
22             std::unordered_set<int> const& finalStates
23         );
24
25         [[nodiscard]] int getInitialState() const;
26
27         [[nodiscard]] bool isValidState(int state) const;
```

```

28     [[nodiscard]] bool isValidCategory(int category) const;
29
30     [[nodiscard]] bool isFinalState(int state) const;
31
32     [[nodiscard]] int getTransition(
33         int state,
34         std::vector<int> const& categories
35     ) const;
36
37 private:
38     const size_t mNoOfStates;           // Q
39     const size_t mNoOfCategories;       // Sigma
40     const std::vector<std::vector<int>>
41         mTransitionTable;                // delta
42     const int mInitialState;             // q_0
43     const std::unordered_set<int> mFinalStates; // F
44 };

```

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 Σ :

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

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.

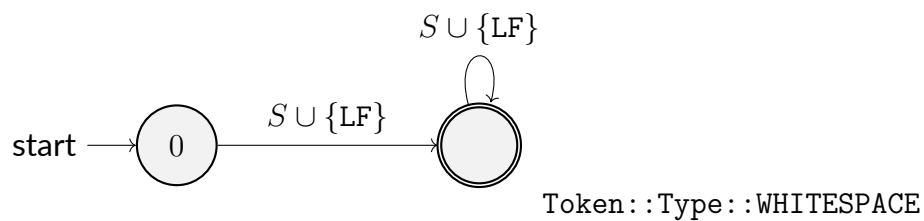


Figure 1: States & transitions for recognising whitespace

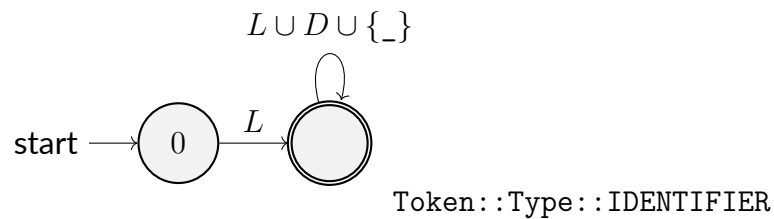


Figure 2: States & transitions for recognising identifiers/keywords

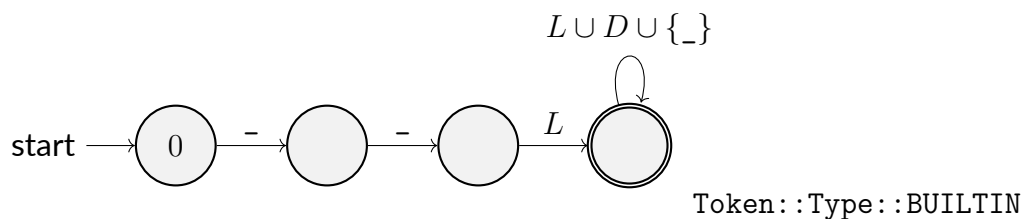


Figure 3: States & transitions for recognising builtins

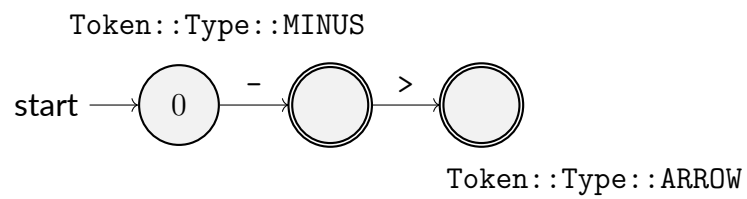


Figure 4: States & transitions for recognising minus and arrow (->)

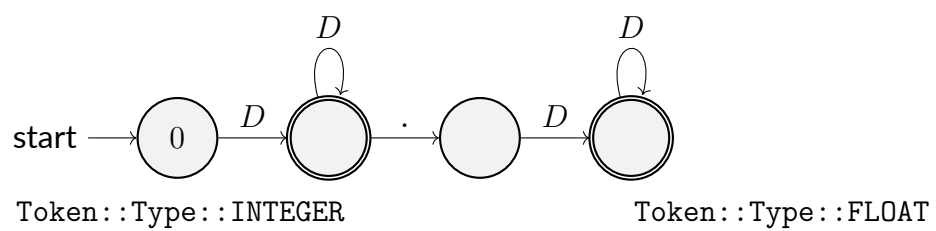


Figure 5: States & transitions for recognising integers and floats

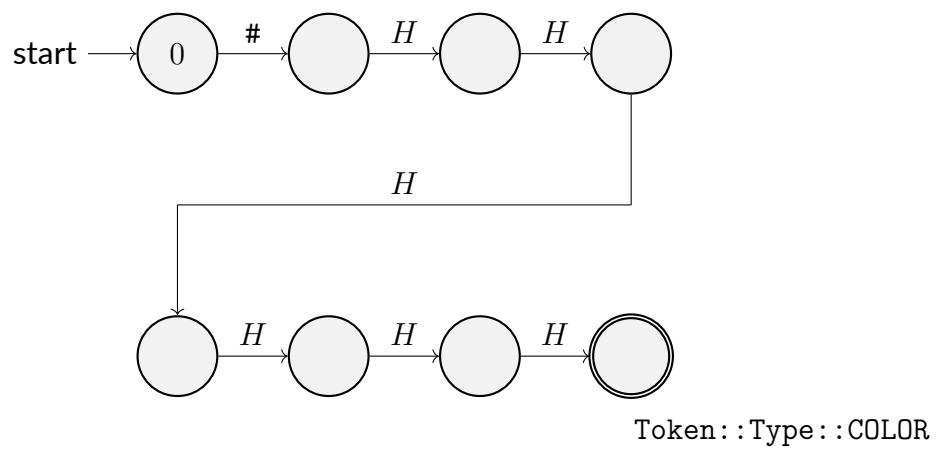


Figure 6: States & transitions for recognising colours

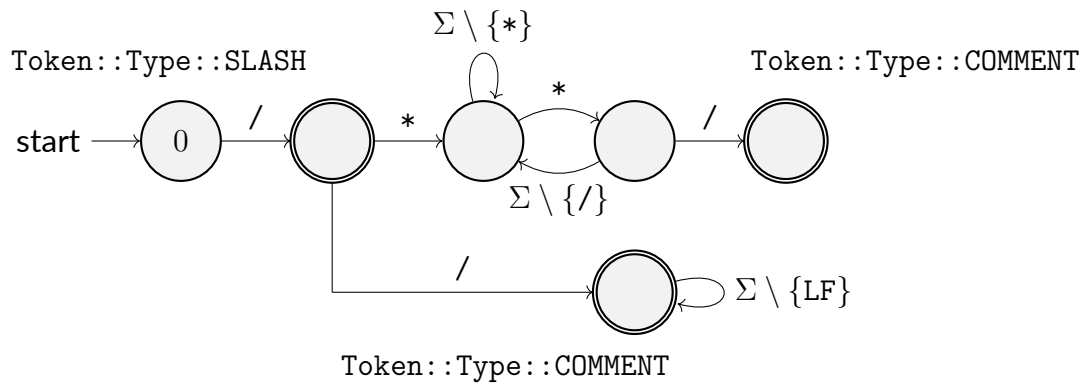


Figure 7: States & transitions for recognising slashes and comments

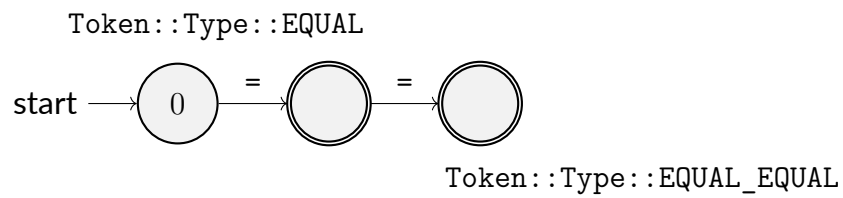


Figure 8: States & transitions for assign and is equal to

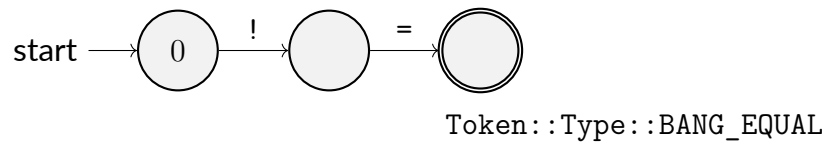


Figure 9: States & transitions for not equal to

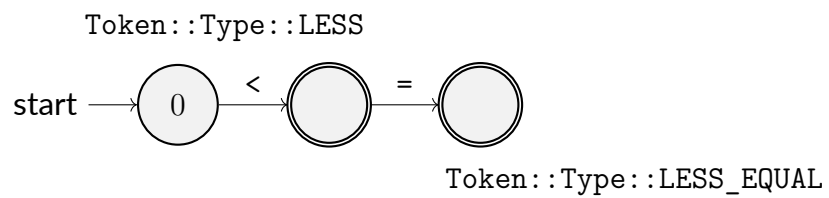


Figure 10: States & transitions for less than and less than or equal to

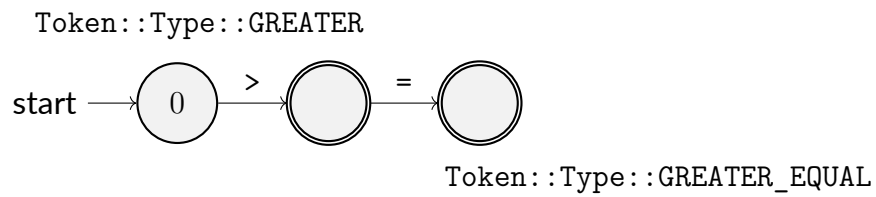


Figure 11: States & transitions for greater than and greater than or equal to

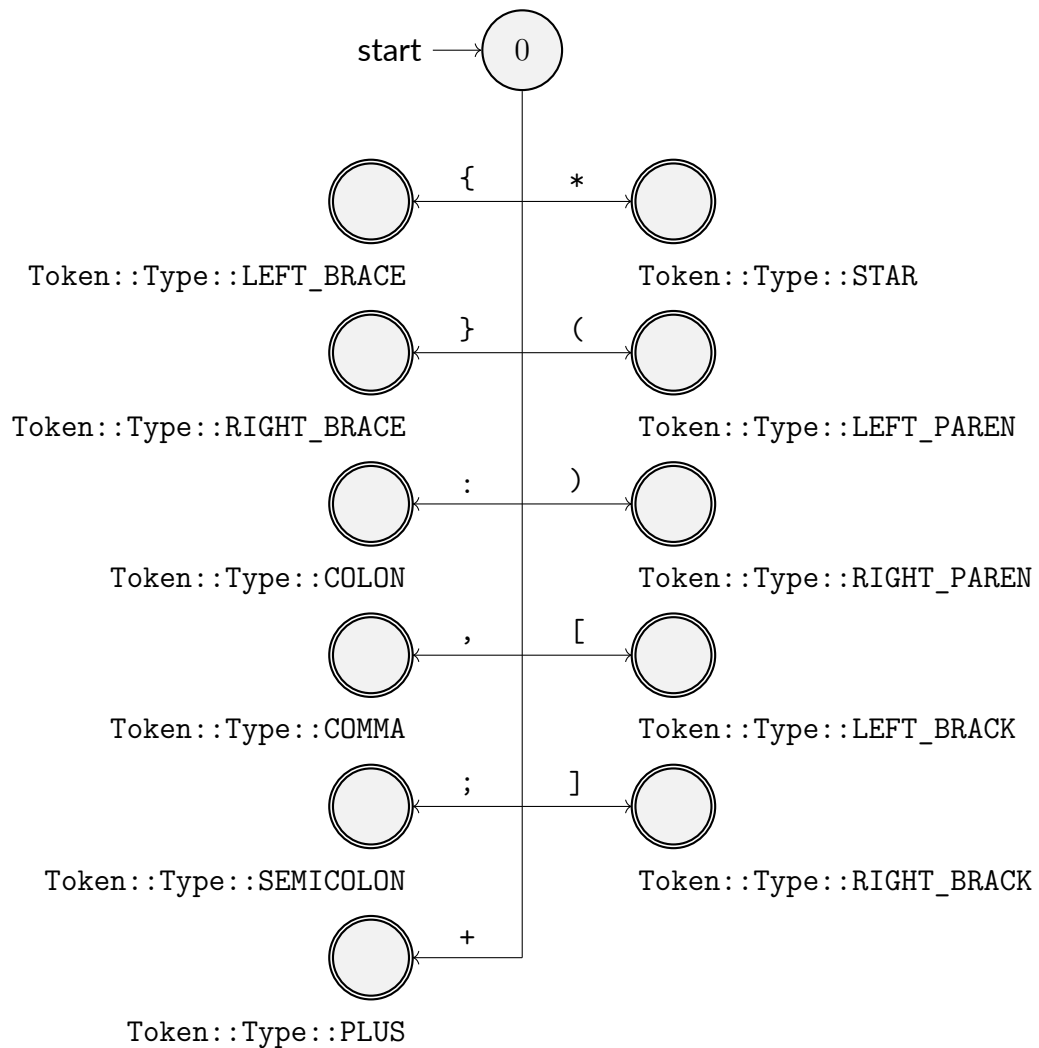


Figure 12: States & transitions for single letter tokens

The Builder & Director

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

```

297 // "/", "//", "/* ... */"
298 builder.addTransition(0, SLASH, 34)
299     .setStateAsFinal(34, Token::Type::SLASH)
300     .addTransition(34, SLASH, 35)
301     .addComplementaryTransition(35, LINEFEED, 35)
302     .setStateAsFinal(35, Token::Type::COMMENT)
303     .addTransition(34, STAR, 36)
304     .addComplementaryTransition(36, STAR, 36)
305     .addTransition(36, STAR, 37)
306     .addComplementaryTransition(37, SLASH, 36)
307     .addTransition(37, SLASH, 38)
308     .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.

```

51     .addCategory(
52         HEX,
53         [](char c) -> bool {
54             return ('0' <= c && c <= '9') ||
55                 ('A' <= c && c <= 'F') ||
56                 ('a' <= c && c <= 'f');
57         }
58     )

```

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.

```
210     // create dfsa
211     Dfsa dfsa(
212         noOfStates,
213         noOfCategories,
214         transitionTable,
215         initialStateIndex,
216         finalStateIndices
217     );
218
219     // create lexer
220     Lexer lexer(
221         std::move(dfsa),
222         std::move(categoryIndexToChecker),
223         std::move(finalStateIndexToTokenType)
224     );
```

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.

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

```
111     for (char ch : lexeme) {
112         mCursor++;
113
114         if (ch == '\n') {
115             mLine++;
116
117             mColumn = 1;
118         } else {
119             mColumn++;
120         }
121     }
122 }
```

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.

```
56     if (state == INVALID_STATE) {
57         mHasError = true;
58
59         fmt::println(
60             stderr,
61             "lexical error at {}:{}:: unexpected "
62             "lexeme '{}'",
63             mLine,
64             mColumn,
65             lexeme
66         );
67     } else {
68         try {
69             token = createToken(
70                 lexeme,
71                 mFinalStateToTokenType.at(state)
72             );
73         } catch (UndefinedBuiltin& error) {
74             mHasError = true;
```

```
75
76         fmt::println(
77             stderr,
78             "lexical error at {}:{}:: {}",
79             mLine,
80             mColumn,
81             error.what()
82         );
83     }
84 }
```

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.

```
22 Runner::Runner(bool dfsaDbg, bool lexerDbg, bool parserDbg)
23     : mDfsaDbg(dfsaDbg),
24       mLexerDbg(lexerDbg),
25       mParserDbg(parserDbg),
26       mLexer(LexerDirector::buildLexer()),
27       mParser(Parser(mLexer)) {
28 }
```

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

2 | The AST & Parsing

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') $\langle \text{Unary} \rangle$ $\langle \text{Primary} \rangle$
$\langle \text{Primary} \rangle$::= $\langle \text{Literal} \rangle$ $\langle \text{SubExpr} \rangle$ $\langle \text{Variable} \rangle$ $\langle \text{ArrayAccess} \rangle$ $\langle \text{FunctionCall} \rangle$
$\langle \text{Program} \rangle$::= { $\langle \text{Stmt} \rangle$ }
$\langle \text{Stmt} \rangle$::= $\langle \text{Block} \rangle$ $\langle \text{VariableDecl} \rangle$ ';' $\langle \text{FunctionDecl} \rangle$ $\langle \text{Assignment} \rangle$ ';' $\langle \text{PrintStmt} \rangle$ ';' $\langle \text{DelayStmt} \rangle$ ';' $\langle \text{WriteBoxStmt} \rangle$ ';' $\langle \text{WriteStmt} \rangle$ ';' $\langle \text{ClearStmt} \rangle$ ';' $\langle \text{IfStmt} \rangle$ $\langle \text{ForStmt} \rangle$

	$\langle \text{WhileStmt} \rangle$ $\langle \text{ReturnStmt} \rangle$ ';'
$\langle \text{Block} \rangle$::= '{' { $\langle \text{Stmt} \rangle$ } '}'
$\langle \text{VariableDecl} \rangle$::= '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 and or 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.

A Better Type System

3 | 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