

# TOC Assignment 2

Akshat Oke                      Aaditya Rathi                      Sankalp Kulkarni  
2020A7PS0284H                      2020A7PS2191H                      2020A7PS1097H

Shreyas Dixit                      Chirag Gadia  
2020A7PS2079H                      2020A7PS1721H

December 2022

## Contents

<b>1</b>	<b>Grammar for Basic C</b>	<b>2</b>
1.1	Grammar syntax . . . . .	2
1.2	The Grammar . . . . .	2
1.3	Comments on the grammar . . . . .	4
1.3.1	Left Associativity . . . . .	4
1.3.2	Expression or Assignment . . . . .	4
<b>2</b>	<b>About the implementation</b>	<b>5</b>
2.1	Types . . . . .	5
2.2	Uninitialized variables . . . . .	5
2.3	Ghost nodes . . . . .	5
<b>3</b>	<b>Remarks</b>	<b>5</b>

# 1 Grammar for Basic C

## 1.1 Grammar syntax

The grammar is mostly in BNF, with the following qualifiers that may follow a grouping or a terminal/non-terminal.

+	One or more of target
*	Zero or more of target
?	Zero or one of target

A grouping of symbols is enclosed by brackets '[' and ']' or '(' and ')'. Additionally, '-' represents negation, which shall *not* match the symbol following it.

Although the qualifiers mentioned above can be converted to BNF fairly easily, negation is problematic and lengthy. For example, to represent -'as' <word>, that is, all words except 'as', we can write

```
<Word-not-as> ::= <notA> <letter>*
| <notS> <letter>*
| 'a' 's' <letter>+
```

```
<notA> ::= [b..z]
```

```
<notS> ::= [a..r] | [t..z]
```

```
<letter> ::= [a..z]
```

[0..9] represents the range of characters 0,1,2,3,4,5,6,7,8,9, similarly for [A..Z] and [a..z].

## 1.2 The Grammar

Start symbol is 'MainProgram'

```
<MainProgram> ::= <Declarations>
| <Program>
| <Declarations> <Program>
```

```
<Declarations> ::= 'int' <identifier> [',' <identifier>]* ';' ;
```

```
<Program> ::= <statement>+
```

```
<statement> ::= ( <assignment>
| <read-statement>
| <write-statement>
```

$$\begin{array}{l} | \langle \text{for-statement} \rangle \\ ) \text{ ' ; ' } \end{array}$$

$$\langle \text{assignment} \rangle ::= \langle \text{identifier} \rangle \text{ '=' } \langle \text{expression} \rangle$$

$$\begin{array}{l} \langle \text{expr-or-assign} \rangle ::= \langle \text{expression} \rangle \\ | \text{ - } \langle \text{binary-expr-strict} \rangle \langle \text{expression} \rangle \text{ '=' } \langle \text{expression} \rangle \end{array}$$

$$\begin{array}{l} \langle \text{binary-expr-strict} \rangle ::= \text{ ' ( ' } \langle \text{binary-expr-strict} \rangle \text{ ' ) ' } \\ | \langle \text{primary-id} \rangle \langle \text{op} \rangle \langle \text{binary-expr} \rangle \end{array}$$

$$\langle \text{binary-expr} \rangle ::= \langle \text{primary-id} \rangle [ \langle \text{op} \rangle \langle \text{binary-expr-strict} \rangle ]?$$

$$\langle \text{primary-id} \rangle ::= \langle \text{identifier} \rangle | \text{ ' ( ' } \langle \text{binary-expr} \rangle \text{ ' ) ' } | \langle \text{integer} \rangle$$

$$\langle \text{op} \rangle ::= \text{ ' + ' } | \text{ ' - ' } | \text{ ' == ' } | \text{ ' * ' } | \text{ ' / ' } | \text{ ' > ' }$$

$$\langle \text{expression} \rangle ::= \langle \text{equality} \rangle$$

$$\langle \text{equality} \rangle ::= \langle \text{comparison} \rangle [ \text{ ' == ' } \langle \text{comparison} \rangle ]?$$

$$\langle \text{comparison} \rangle ::= \langle \text{term} \rangle [ \text{ ' > ' } \langle \text{term} \rangle ]?$$

$$\langle \text{term} \rangle ::= \langle \text{factor} \rangle [ ( \text{ ' + ' } | \text{ ' - ' } ) \langle \text{factor} \rangle ]^*$$

$$\langle \text{factor} \rangle ::= \langle \text{unary} \rangle [ ( \text{ ' * ' } | \text{ ' / ' } ) \langle \text{unary} \rangle ]^*$$

$$\langle \text{unary} \rangle ::= \langle \text{primary} \rangle$$

$$\langle \text{primary} \rangle ::= \langle \text{integer} \rangle | \langle \text{identifier} \rangle | \text{ ' ( ' } \langle \text{expression} \rangle \text{ ' ) ' }$$

$$\langle \text{read-statement} \rangle ::= \text{ ' read ' } \langle \text{identifier} \rangle$$

$$\langle \text{write-statement} \rangle ::= \text{ ' write ' } ( \langle \text{integer} \rangle | \langle \text{identifier} \rangle )$$

$$\langle \text{for-statement} \rangle ::= \text{ ' for ' } \text{ ' ( ' } \langle \text{expr-or-assign} \rangle ? \text{ ' ; ' } \langle \text{expression} \rangle ? \text{ ' ; ' } \langle \text{expr-or-assign} \rangle ? \text{ ' ) ' } \langle \text{block} \rangle$$

$$\begin{aligned}
\langle block \rangle &::= \{ \langle statement \rangle^* \} \\
\langle identifier \rangle &::= - \langle keyword \rangle \langle alpha \rangle + \\
\langle integer \rangle &::= \langle digit \rangle + \\
\langle digit \rangle &::= [0..9] \\
\langle alpha \rangle &::= [a..z] \\
\langle keyword \rangle &::= \text{'for'} \mid \text{'int'} \mid \text{'read'} \mid \text{'write'}
\end{aligned}$$

### 1.3 Comments on the grammar

It is possible to skip the  $\langle equality \rangle$  rule as it only forwards  $\langle expression \rangle$ , but makes the meaning clearer.

#### 1.3.1 Left Associativity

The grammar is not left-associative as required, for it would require it to be left-recursive as well. Since the parsing algorithm used is recursive descent, the additive and multiplicative productions are flattened out and interpreted left-recursively by the interpreter.

#### 1.3.2 Expression or Assignment

The 'for' clauses 1 and 3 were specified to be either expressions or assignment statements. In most programming languages today, the assignment is part of the expression, with the lowest priority, whose return value is the expression on its RHS. This in grammar, it looks like

$$\begin{aligned}
\langle expression \rangle &::= \langle assignment-expression \rangle \\
\langle assignment-expression \rangle &::= \langle equality \rangle \\
&\quad \mid \langle identifier \rangle '=' \langle assignment-expression \rangle
\end{aligned}$$

The problem with this due to the fact that  $\langle equality \rangle$  may also begin with an  $\langle identifier \rangle$  rule, failing the pairwise-disjointness test. This may be then left-factored to allow LL(1) parsers to succeed. (Otherwise a lookahead of 2 tokens is required)

Since expressions cannot have assignments, a new rule  $\langle expr\text{-or-assign} \rangle$  had to be made to rule out ambiguity in choosing which branch to take. If the LHS  $\langle expression \rangle$  was  $\langle identifier \rangle$  as required by the specification, it would not be possible to parse it as an expression (once '=' was matched)

without backtracking. Thus, it an parses expression, looks for a '=', and if present, attempts to negate <binary-expr-strict> (not match binary expressions involving identifiers) and then proceed to parse the RHS. (This rule prevents the grammar from generating expressions for assignment targets)

Just replacing the Kleene star by '+' in <comparison> and <equality> would have resulted in Python-style expressions such as `a>b>c!`

## 2 About the implementation

The parser is an LL(1) recursive descent parser.

### 2.1 Types

Since the document mentions nothing about types, this implementation assumes 'Basic C' to be strongly typed, with two types.

There are two types of expressions: number and boolean. Only a boolean expression can be used in the condition clause of the for statement.

Printing a boolean variable will print 'true' or 'false' correspondingly. Booleans cannot be operated upon by arithmetic operations.

### 2.2 Uninitialized variables

The <Declarations> rule declares variables that can be used in the program. Variables have been implemented with a hashmap. Referencing a variable before assigning a value to it causes an 'Uninitialized variable' runtime error.

It was possible to implicitly give variables an initial value of 0, but this may lead to unclear code.

### 2.3 Ghost nodes

The for statement grammar may not include all three clauses, for example `for(;a<4;)`. In such cases the parser inserts a 'ghost node' which evaluates to 'true' but does not get printed by the AST printer code.

This ghost node is also used in <expr-or-assign> in order to provide easy access to the lvalue identifier lexeme.

## 3 Remarks

Before the correction in the for statement syntax, all three clauses were expressions and had no side effects. This means they could be executed any number of times and in any order during the loop execution. Although it seemed that we could get away easily by not using the AST for interpreting, the additional "rewind" functionality required in the scanner was better left alone.