

CSE423 Assignment 1 Documentation

Rebecca Castillo
Geoff Knox
Krishna Marentes
Elijah Orozco

February 16 2020

1 Program Overview

This file serves as the documentation for our compiler, currently in the scanning and parsing stages. Our user guide can be found in “usage.pdf”. Our compiler is written in Java and currently consists of three parts:

- Grammar
- Scanner/Parser Generator
- Driver

A third-party tool, “Antlr,” serves as the scanner/parser generator that takes our Grammar (stored in the .g4 file) as input. Antlr is prompted to generate and run scanning and parsing files with the specified grammar from the Driver file, which handles command line arguments, the compiler/Antlr interface, and general “main” function duties.

2 Design - Antlr

Our compiler uses Antlr4, a tool similar to Flex/Bison that can be used for compilers written in Java. Antlr is both a scanner and parser all at once, and uses a simple format for specifying the language to parse.

We originally began our design handwriting our scanner code, with the same plan for the parser. However, Antlr was an attractive choice for us due to the advantages outlined below.

2.1 Advantages

Antlr allowed us to fast-track our parser development by greatly reducing the amount of code we needed to design and write. Antlr handles the parsing algorithm so that we can avoid writing tedious and error-prone code ourselves.

Additionally, the grammar files Antlr uses are simple and easy to write. There is also a convenient plugin for our team's IDE of choice, IntelliJ, that seamlessly incorporates Antlr into our development environment, including a useful dynamic parse tree display.

2.2 Disadvantages

As simple as Antlr grammar files are and as much function as Antlr provides, there is still the learning curve associated with using and integrating such a complex tool. We are also aware that Antlr grammar files can be difficult to debug, but we are prepared to handle such issues, especially with the tools the IntelliJ Antlr plugin provides.

3 Language Specification

3.1 Introduction

Here are the types of the various elements by type font or symbol for the grammar that follows:

- **Terminals, including terminal punctuation, are bolded**
- **TOKENS ARE ALL CAPS**
- Nonterminals are in “normal” font
- The symbol ϵ means the empty string
- “EOF” indicates the end of the file

The grammar also uses standard regex terms. Symbols indicating regex rules are never bolded so as not to be confused with terminal symbols. Parentheses are also used to avoid confusion. For example, an instance of “(*)” indicates that the * terminal may be seen zero or more times.

3.2 Tokens

- $ID \rightarrow (- \mid [a-zA-Z])^+ ([a-zA-Z] \mid DIGIT \mid _)^*$
- $CHARCONST \rightarrow ' CHARCHARS^+ '$
- $STRINGCONST \rightarrow " STRINGCHARS^* "$
- $INT \rightarrow DIGIT^+ \mid (0x \mid 0X)HEXDIGIT^+ \mid 0OCTALDIGIT^+ \mid (0b \mid 0B)BINARYDIGIT^+ \mid FLOAT$

The following are “helper tokens” that are not strictly considered tokens by our compiler but serve to help define the tokens above.

- CHARCHARS \rightarrow
 $(\sim [\backslash'\\r\\n] \mid '\backslash' (\cdot \mid \text{EOF}))$
- LETTER $\rightarrow [\text{a-zA-Z}]$
- DIGIT $\rightarrow [\text{0-9}]$
- HEXDIGIT $\rightarrow [\text{0-9A-Fa-f}]$
- OCTALDIGIT $\rightarrow [\text{0-7}]$
- BINARYDIGIT $\rightarrow [\text{0-1}]$
- FLOAT $\rightarrow [\text{0-9}]^+ \cdot [\text{0-9}]^+ \text{EXP}?(f|F)?$
 $\mid \cdot [\text{0-9}]^+ \text{EXP}?(f|F)?$
 $\mid [\text{0-9}]^+ \text{EXP}(f|F)?$
 $\mid [\text{0-9}]^+ (f|F)?$
- EXP $\rightarrow (e|E) (+|-)? [\text{0-9}]^+$

3.3 The Grammar

1. program \rightarrow declarationList **EOF**
2. declarationList \rightarrow declarationList declaration \mid declaration
3. declaration \rightarrow varDeclaration \mid funDeclaration \mid structDeclaration \mid enumDeclaration \mid ;
4. structDeclaration \rightarrow **static?** **struct** ID { unInitVarDecl* } ID? ;
5. structInit \rightarrow **static?** **struct** ID (*)* varDeclList ;
6. enumDeclaration \rightarrow **enum** ID { enumDeclList } ID? ;
7. enumDeclList \rightarrow enumDeclList , enumId \mid enumId \mid ϵ
8. enumId \rightarrow ID ASSIGNMENT enumExpression \mid ID
9. enumInit \rightarrow **enum** ID ID ; \mid **enum** ID ID = expression ;
10. unInitVarDecl \rightarrow typeSpecifier unInitVarDeclList ;
11. unInitVarDeclList \rightarrow unInitVarDeclList , varDeclId \mid varDeclId \mid ϵ
12. varDeclaration \rightarrow typeSpecifier varDeclList ; \mid scopedVarDeclaration ; \mid structInit \mid enumInit

13. $\text{scopedVarDeclaration} \rightarrow \text{scopedTypeSpecifier varDeclList}$
14. $\text{forLoopVars typeSpecifier varDeclList} \mid \text{expression List}$
15. $\text{varDeclList} \rightarrow \text{varDeclList} , \text{varDeclInitialize} \mid \text{varDeclInitialize}$
16. $\text{varDeclInitialize} \rightarrow \text{varDeclId} \mid \text{varDeclId} = (\text{expression} \mid \{ \text{expressionList} \})$
17. $\text{varDeclId} \rightarrow \text{ID} ([\text{expression?}])^*$
18. $\text{scopedTypeSpecifier} \rightarrow \text{static typeSpecifier} \mid \text{typeSpecifier}$
19. $\text{typeSpecifier} \rightarrow (\text{int} \mid \text{float} \mid \text{double} \mid \text{char} \mid \text{long} \mid \text{unsigned} \mid \text{signed} \mid \text{void} \mid \text{short})(*)^*$
20. $\text{funDeclaration} \rightarrow \text{typeSpecifier ID} (\text{params}) (\text{compoundStmt} \mid ;)$
21. $\text{params} \rightarrow \text{params} , \text{parameter} \mid \text{parameter} \mid \epsilon$
22. $\text{parameter} \rightarrow \text{typeSpecifier paramId}$
23. $\text{paramId} \rightarrow \text{ID} ([])^*$
24. $\text{statement} \rightarrow \text{expressionStmt} \mid \text{compoundStmt} \mid \text{selectionStmt} \mid \text{iterationStmt} \mid \text{returnStmt} \mid \text{breakStmt} \mid \text{gotoStmt} \mid \text{labelStmt}$
25. $\text{structExpressionList} \rightarrow \text{expressionList} , .? \text{expression} \mid .? \text{expression} \mid \epsilon$
26. $\text{expressionList} \rightarrow \text{expressionList} , \text{expression} \mid \text{expression} \mid \epsilon$
27. $\text{expressionStmt} \rightarrow \text{expression} ; \mid ;$
28. $\text{compoundStmt} \rightarrow \{ \text{statementList} \}$
29. $\text{statementList} \rightarrow \text{statementList} (\text{statement} \mid \text{varDeclaration}) \mid \mathbf{EPS}$
30. $\text{defaultList} \rightarrow \text{statementList} (\text{statement} \mid \text{varDeclaration}) \mid \text{statement}$
31. $\text{elsifList} \rightarrow \text{elsifList} \mathbf{else if} (\text{expression}) \text{statement} \mid \mathbf{EPS}$
32. $\text{selectionStmt} \rightarrow \mathbf{if} (\text{expression}) \text{statement} \text{elsifList} \mid \mathbf{if} (\text{expression}) \text{statement} \text{elsifList} \mathbf{else} \text{statement} \mid \mathbf{switch} (\text{expression}) \text{switchCase} \text{defaultList} \mid \mathbf{switch} (\text{expression}) \mathbf{default} : \text{defaultList} \mid \mathbf{switch} (\text{expression}) \{ \text{switchList} (\mathbf{default} : \text{defaultList})? \}$
33. $\text{switchList} \rightarrow \text{switchList} \text{switchCase} \mid \text{switchCase} \mid \epsilon$
34. $\text{switchCase} \rightarrow \mathbf{case} (\mathbf{INT} \mid \mathbf{CHARCONST}) : (\text{defaultList} \mid \text{statementList})$
35. $\text{iterationStmt} \rightarrow \mathbf{while} (\text{expression}) \text{statement} \mid \mathbf{do} \text{statement} \mathbf{while} (\text{expression}); \mid \mathbf{for} (\text{forLoopVars} ; \text{expressionList} ; \text{expressionList}) \text{statement}$

- 36. returnStmt \rightarrow **return** ; | **return** expression ;
- 37. breakStmt \rightarrow **break** ;
- 38. gotoStmt \rightarrow **goto** labelId ;
- 39. labelStmt \rightarrow labelId :
- 40. labelId \rightarrow ID
- 41. expression \rightarrow mutable = expression
 - | mutable += expression
 - | mutable -= expression
 - | mutable *= expression
 - | mutable /= expression
 - | mutable %= expression
 - | mutable <<= expression
 - | mutable >>= expression
 - | mutable &= expression
 - | mutable |= expression
 - | mutable ^= expression
 - | (mutable | immutable) << expression
 - | (mutable | immutable) >> expression
 - | (mutable | immutable) & expression
 - | (mutable | immutable) | expression
 - | (mutable | immutable) ^ expression
 - | simpleExpression
- 42. enumExpression \rightarrow properUnaryOps INT << enumExpression
 - | properUnaryOps INT >> enumExpression
 - | properUnaryOps INT & enumExpression
 - | properUnaryOps INT | enumExpression
 - | properUnaryOps INT ^ enumExpression
 - | properUnaryOps INT || enumExpression
 - | properUnaryOps INT && enumExpression
 - | properUnaryOps INT < enumExpression
 - | properUnaryOps INT <= enumExpression
 - | properUnaryOps INT > enumExpression
 - | properUnaryOps INT >= enumExpression
 - | properUnaryOps INT + enumExpression
 - | properUnaryOps INT - enumExpression
 - | properUnaryOps INT * enumExpression
 - | properUnaryOps INT
- 43. properUnaryOps \rightarrow (- | ~ | !)*
- 44. simpleExpression \rightarrow (simpleExpression || andExpression) | andExpression

45. $\text{andExpression} \rightarrow \text{andExpression} \ \&\& \ \text{unaryRelExpression} \mid \text{unaryRelExpression}$
46. $\text{unaryRelExpression} \rightarrow ! \ \text{unaryRelExpression} \mid \text{relExpression}$
47. $\text{relExpression} \rightarrow \text{sumExpression} \ \text{relop} \ \text{sumExpression} \mid \text{relExpression} \ \text{relop} \ \text{relExpression} \mid \text{sumExpression}$
48. $\text{relop} \rightarrow <= \mid < \mid > \mid >= \mid == \mid !=$
49. $\text{sumExpression} \rightarrow \text{sumExpression} \ \text{sumop} \ \text{mulExpression} \mid \text{mulExpression}$
50. $\text{sumop} \rightarrow + \mid -$
51. $\text{mulExpression} \rightarrow \text{mulExpression} \ \text{mulop} \ \text{unaryExpression} \mid \text{unaryExpression}$
52. $\text{mulop} \rightarrow * \mid / \mid \%$
53. $\text{unaryExpression} \rightarrow \text{unaryop} \ \text{unaryExpression} \mid \text{mutable} \ ++ \mid \text{mutable} \ - \mid - \ \text{mutable} \mid ++ \ \text{mutable} \mid \text{factor}$
54. $\text{unaryop} \rightarrow - \mid * \mid ! \mid \& \mid \sim$
55. $\text{factor} \mid \rightarrow \text{immutable} \mid \text{mutable}$
56. $\text{mutable} \rightarrow (*)^* \ \text{ID} \mid \text{mutable} \ [\ \text{expression} \] \mid \text{mutable} \ (\cdot \mid ->) \ \text{mutable} \mid \text{immutable} \ (\cdot \mid ->) \ \text{mutable} \mid (\ \text{expression} \) \ (\cdot \mid ->) \ \text{mutable} \mid (*)^* \ (\ \text{expression} \)$
57. $\text{immutable} \rightarrow (\ \text{expression} \) \mid \text{call} \mid \text{constant}$
58. $\text{call} \rightarrow \text{ID} \ (\ \text{args} \) \mid \text{sizeof} \ (\ (\text{struct} \ \text{ID} \ (*)^* \mid \text{typeSpecifier} \mid \text{ID}) \)$
59. $\text{args} \rightarrow \text{argList} \mid \epsilon$
60. $\text{argList} \rightarrow \text{argList} \ , \ \text{expression} \mid \text{expression}$
61. $\text{constant} \rightarrow \text{INT} \mid \text{CHARCONST} \mid \text{STRINGCONST}$

3.4 Semantic Notes

- HEX, OCTAL, and BINARYDIGIT default to **int** when parsed
- Many variables can be declared and/or initialized in one statement, e.g. “int a = 1, b = 2;”

3.5 Limitations

The following are not supported by our grammar.

- Preprocessor statements
- Casting
- Ternary operations

We've attempted to implement the following, but support can be considered to be in "beta mode" as there may be edge cases that have not been tested.

- Pointers
- Arrays
- Strings