

SCC by KREG

“Not-very-good C Compiler”

Krishna Marentes
Rebecca Castillo
Elijah Orozco
Geoff Knox

May 3, 2020

Contents

1	Program Overview	2
2	Compiler Capabilities and Limitations	2
2.1	Capabilities	2
2.2	Limitations	3
3	Scanning and Parsing - Antlr	3
3.1	Overview	3
3.2	Advantages	4
3.3	Disadvantages	4
4	The Syntax Tree	4
4.1	Overview	4
4.2	Advantages	5
4.3	Disadvantages	5
5	Language Specification	5
5.1	Introduction	5
5.2	Tokens	5
5.3	The Grammar	6
5.4	Semantic Notes	10
5.5	Limitations	10
6	Symbol Table	10
7	Intermediate Language Design	11
7.1	Overview	11

7.2	Limitations	12
8	Optimizations	13
8.1	Basic Blocks	13
8.2	Methods	13
9	Assembly Generation	15
9.1	Calling Convention	15
9.2	Register Allocation	15

1 Program Overview

This file serves as the documentation for our compiler, which currently can scan and parse C code and generate a linear IR. Our user guide can be found in “usage.pdf”. Our compiler is written in Java and consists of three parts:

- Phase 1: Scanning/Parsing with C Grammar
 - Input: C code
 - Output: Abstract Syntax Tree
- Phase 2: IR Generation
 - Input: Abstract Syntax Tree
 - Output: Linear intermediate representation
- Phase 3: Optimization and Assembly Generation
 - Input: Linear intermediate representation
 - Output: Optimized IR, x86 AT&T assembly

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 Compiler Capabilities and Limitations

2.1 Capabilities

The following list states everything that our compiler supports

- Identifiers, variables, functions
- Keywords
- Arithmetic expressions

- Assignment
- Boolean expressions (excluding `||` and `&&`)
- Goto statement
- If / else control flow (nesting works also)
- Unary operators (`-`, `~`)
- Return statements
- While loops
- `++`, `-`, `-=`, `+=`, etc
- Binary operators and assignment (`<<`, `<<=`, `|`, `|=`, `&` etc)

2.2 Limitations

The following states what our compiler can not do

- Types other than integers. Originally this was supported during scanning, parsing, and IR generation, but it was not something we could convert into assembly with the time we had remaining.
- For loops. This was supported during scanning and parsing, but dropped in subsequent projects.
- Switch statements. Supported during scanning and parsing, dropped in subsequent projects.
- Pointers, arrays, strings. Somewhat implemented in scanning and parsing, no longer pursued in project 2 and later.
- Struct, enum. Same as above.
- Casting, type promotion. Never implemented in any project.
- Multiple boolean expressions not allowed in if/else blocks or while loops. For example, `if(i < 7 || f > 2)` is not allowed. However, `if(i > 7)` is completely fine.

3 Scanning and Parsing - Antlr

3.1 Overview

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, 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.

3.2 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.

3.3 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.

It was discovered when beginning Phase 2 that Antlr left us with no convenient way to convert the automatically-generated parse tree into an abstract syntax tree. Therefore, some extra effort was required to convert Antlr's tree data structure into our own design for AST generation and conversion to IR.

4 The Syntax Tree

4.1 Overview

As discussed before, Antlr generates a thorough parse tree of the source code according to the defined grammar, but leaves no obvious way to use the generated parse tree. Therefore, our compiler defined its own tree data structure and methods to convert the parse tree (defined in Antlr's own tree data structure) into our custom tree design. This conversion also trims the parse tree into an abstract syntax tree (AST) for clarity and convenience for generating the IR later. Tree conversion routines are held in the `ASTNode` class.

To aid IR generation, each AST node is resolved into a certain subclass node according to the grammar. For example, a node holding the "funDeclaration" grammar rule in the parse tree would be resolved into a node of type "FunDeclaration" in the AST. The AST uses seven subclasses of AST nodes: Program, TypeSpecifier, VarDeclaration, FunDeclaration, StructDeclaration, Expression, and Statement. (Note that the StructDeclaration is not fully implemented in the current version.) Each subclass is a customized node with relevant members (such as "params" for FunDeclaration) that in some cases reduces the total number of nodes in the AST.

4.2 Advantages

The advantages of this design center on the ease of AST printing and IR generation. IR generation routines are customized for each subclass, so that IR generation can be successfully performed by recursively walking the tree and checking for what instance each node belongs to. This design is also useful when using an IDE debugger, since each node will have relevant members directly in the object, rather than needing to look through layers of child and parent nodes.

4.3 Disadvantages

For all the advantages, disadvantages do exist for our approach. For one, the overall design is somewhat complex. Using the AST requires fairly deep knowledge of each subclass as well as the grammar itself. In addition, it can be confusing to know when to look for information in a node's members or in its children. However, we've found the current design to be sufficiently helpful for IR

5 Language Specification

5.1 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.

Single-line comments (denoted by “//”) and block comments (denoted by “/* */”) are successfully ignored by the grammar.

5.2 Tokens

- $ID \rightarrow (- \mid [a-zA-Z])^+ ([a-zA-Z] \mid DIGIT \mid -)^*$
- $CHARCONST \rightarrow ' CHARCHARS^+ '$

- $\text{STRINGCONST} \rightarrow " \text{STRINGCHARS}^* "$
- $\text{INT} \rightarrow \text{DIGIT}^+$
 $\quad | (\text{0x} \mid \text{0X})\text{HEXDIGIT}^+$
 $\quad | \text{OCTALDIGIT}^+$
 $\quad | (\text{0b} \mid \text{0B})\text{BINARYDIGIT}^+$
 $\quad | \text{FLOAT}$

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

- $\text{CHARCHARS} \rightarrow$

$$(\sim [\backslash'\\r\n] \mid '\backslash' (. \mid \text{EOF}))$$

- $\text{LETTER} \rightarrow [\text{a-zA-Z}]$
- $\text{DIGIT} \rightarrow [\text{0-9}]$
- $\text{HEXDIGIT} \rightarrow [\text{0-9A-Fa-f}]$
- $\text{OCTALDIGIT} \rightarrow [\text{0-7}]$
- $\text{BINARYDIGIT} \rightarrow [\text{0-1}]$
- $\text{FLOAT} \rightarrow [\text{0-9}]^+ . [\text{0-9}]^+ \text{EXP}?(f|\text{F})?$
 $\quad | . [\text{0-9}]^+ \text{EXP}?(f|\text{F})?$
 $\quad | [\text{0-9}]^+ \text{EXP}(f|\text{F})?$
 $\quad | [\text{0-9}]^+ (f|\text{F})?$
- $\text{EXP} \rightarrow (\text{e}|\text{E}) (+|-)? [\text{0-9}]^+$

5.3 The Grammar

1. $\text{program} \rightarrow \text{declarationList } \text{EOF}$
2. $\text{declarationList} \rightarrow \text{declarationList } \text{declaration} \mid \text{declaration}$
3. $\text{declaration} \rightarrow \text{varDeclaration} \mid \text{funDeclaration} \mid \text{structDeclaration} \mid \text{enumDeclaration} \mid ;$
4. $\text{structDeclaration} \rightarrow \text{static? } \text{struct } \text{ID } \{ \text{unInitVarDecl}^* \} \text{ID? ;}$
5. $\text{structInit} \rightarrow \text{static? } \text{struct } \text{ID } (*)^* \text{varDeclList ;}$
6. $\text{enumDeclaration} \rightarrow \text{enum } \text{ID } \{ \text{enumDeclList} \} \text{ID? ;}$
7. $\text{enumDeclList} \rightarrow \text{enumDeclList } , \text{enumId} \mid \text{enumId} \mid \epsilon$

8. $\text{enumId} \rightarrow \text{ID ASSIGNMENT enumExpression} \mid \text{ID}$
9. $\text{enumInit} \rightarrow \mathbf{enum} \text{ ID ID } ; \mid \mathbf{enum} \text{ ID ID } = \text{expression} ;$
10. $\text{unInitVarDecl} \rightarrow \text{typeSpecifier unInitVarDeclList} ;$
11. $\text{unInitVarDeclList} \rightarrow \text{unInitVarDeclList} , \text{varDeclId} \mid \text{varDeclId} \mid \epsilon$
12. $\text{varDeclaration} \rightarrow \text{typeSpecifier varDeclList} ; \mid \text{scopedVarDeclaration} ; \mid \text{structInit} \mid \text{enumInit}$
13. $\text{scopedVarDeclaration} \rightarrow \text{scopedTypeSpecifier varDeclList}$
14. $\text{forLoopVars} \text{ typeSpecifier varDeclList} \mid \text{expression List}$
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. $\text{scopedTypeSpecifier} \rightarrow \mathbf{static} \text{ typeSpecifier} \mid \text{typeSpecifier}$
19. $\text{typeSpecifier} \rightarrow (\mathbf{int} \mid \mathbf{float} \mid \mathbf{double} \mid \mathbf{char} \mid \mathbf{long} \mid \mathbf{unsigned} \mid \mathbf{signed} \mid \mathbf{void} \mid \mathbf{short})(^*)^*$
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. $\text{statementList} \rightarrow \text{statementList} (\text{statement} \mid \text{varDeclcaration}) \mid \mathbf{EPS}$
30. $\text{defaultList} \rightarrow \text{statementList} (\text{statement} \mid \text{varDeclaration}) \mid \text{statement}$
31. $\text{elseifList} \rightarrow \text{elseifList} \mathbf{else if} (\text{expression}) \text{statement} \mid \mathbf{EPS}$

32. selectionStmt \rightarrow **if** (expression) statement elsifList
| **if** (expression) statement elsifList **else** statement
| **switch** (expression) switchCase defaultList
| **switch** (expression) **default** : defaultList
| **switch** (expression) { switchList (**default** : defaultList)? }
33. switchList \rightarrow switchList switchCase | switchCase | ϵ
34. switchCase \rightarrow **case** (**INT** | **CHARCONST**) : (defaultList | statementList)
35. iterationStmt \rightarrow **while** (expression) statement
| **do** statement **while** (expression);
| **for** (forLoopVars ; expressionList ; expressionList) 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. andExpression \rightarrow andExpression && unaryRelExpression | unaryRelExpression
- 46. unaryRelExpression \rightarrow ! unaryRelExpression | relExpression
- 47. relExpression \rightarrow sumExpression relop sumExpression | relExpression relop relExpression | sumExpression
- 48. relop \rightarrow <= | < | > | >= | == | !=
- 49. sumExpression \rightarrow sumExpression sumop mulExpression | mulExpression
- 50. sumop \rightarrow + | -
- 51. mulExpression \rightarrow mulExpression mulop unaryExpression | unaryExpression
- 52. mulop \rightarrow * | / | %
- 53. unaryExpression \rightarrow unaryop unaryExpression | mutable ++ | mutable - | - mutable | ++ mutable | factor
- 54. unaryop \rightarrow - | * | ! | & | ~
- 55. factor | \rightarrow immutable | mutable
- 56. mutable \rightarrow (*)* ID
 - | mutable [expression]
 - | mutable (. | ->) mutable
 - | immutable (. | ->) mutable
 - | (expression) (. | ->) mutable
 - | (*)* (expression)
- 57. immutable \rightarrow (expression) | call | constant
- 58. call \rightarrow ID (args) | sizeof ((struct ID (*)* | typeSpecifier | ID))

- 59. $\text{args} \rightarrow \text{argList} \mid \epsilon$
- 60. $\text{argList} \rightarrow \text{argList} , \text{expression} \mid \text{expression}$
- 61. $\text{constant} \rightarrow \text{INT} \mid \text{CHARCONST} \mid \text{STRINGCONST}$

5.4 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;”

5.5 Limitations

The following are not supported by our grammar.

- Preprocessor statements
- Casting
- Ternary operations

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

- Pointers
- Arrays
- Strings

6 Symbol Table

Our compiler utilizes a symbol table to keep track of declared variables and functions and their types. The symbol table is built from the complete abstract syntax tree. The general structure of our symbol table with examples is shown:

Name	Type
i	int
c	char
main	int

The “Name” field of our symbol table can be of two different types: a plain **SymbolEntry** type for variable declarations, or a **SymbolTable** type for function declarations. This way, function return types are stored in the symbol table belonging to the scope they are declared in, and each function’s scoped variables can be found by looking inside the symbol table of the function’s name.

For example, the following C code

```

int global_var;
char foo() {
    char i;
    return i;
}
int main() {
    int i;
    char j;
    return i;
}

```

produces the following symbol table:

Name		Type
global_var		int
foo		char
Name	Type	
i	char	
main		int
Name	Type	
i	int	
j	char	

The symbol table mainly serves to check for type and declaration errors in the source code. For example, if the source code tries to reference variable “i” and “i” is not found in the symbol table, the compiler will throw an error to the user. It should be noted that error checking is in early stages—this feature has not been thoroughly tested.

7 Intermediate Language Design

7.1 Overview

Our linear intermediate representation (IR) follows a simple format not far from the original C code. It flattens all loops with `goto` statements, and uses `if` and `else` to handle conditional jumps. Our IR nearly follows the single static assignment (SSA) form.

Temporary variables and labels are named with the prefix “KREG” followed by a period and unique integer: “KREG.0”. Labels are differentiated from variables by angle brackets: “<KREG.0>”. We rely on the period in our temporary variable naming scheme to avoid naming collisions between IR-generated variables and C code variables.

Statements end with a semicolon, and labels end with a colon. Function bodies are denoted with curly brackets, each on their own line. Function parameters are denoted in parenthesis next to the function name, similar to C. All function declarations are preceded by the keyword “function”. Though our algorithm does sometimes result

in seemingly redundant temporary variable assignments (see example below) by using pseudo-SSA, we find that this method smoothly and easily breaks the original code down into manageable lines for optimization.

For the following line of C code:

```
i++;
```

our compiler will produce the following IR:

```
KREG.1 = i;  
i = i + 1;  
KREG.2 = KREG.1;
```

Although in this case our IR unnecessarily inflates the original C code, it has no real impact on the program's function. More importantly, producing the IR in this way allows us to use the proper values of variables for unary expressions. For example, if the C code was changed to

```
x = i++;
```

then our IR would produce

```
KREG.1 = i;  
i = i + 1;  
KREG.2 = KREG.1;  
x = KREG.2;
```

which correctly assigns the un-incremented value of `i` to `x`.

Currently, the IR is stored as a single string, which each line delineated by the running system's *End of Line* marker.

7.2 Limitations

The following features, in addition to those not supported by the grammar, are not supported by our IR:

- For loops
- Pointers
- Strings
- Structures and Enums
- Arrays
- Multiple boolean expressions in `if` statements and `while` loops. E.g.

```
while (i < 10 && j > 10)
```

Currently, our compiler only supports two levels of symbol scoping: the global scope, and scopes for functions. Therefore, loops do *not* have their own scope.

8 Optimizations

Our compiler performs the following optimizations on the linear IR: constant folding, constant propagation, algebraic simplification, and identity removal. In addition, removal of "dead code" is mostly implemented, but technical issues prevent its use (see Details below.)

8.1 Basic Blocks

The SCC Optimizer relies on isolating basic blocks of execution from the IR. The instructions in a basic block can be optimized together, but each block is examined separately; no code outside the block can impact optimizations.

Basic blocks are identified as blocks with only one entry point and one exit point. These points are defined in SCC as follows:

Entry points: Functions, labels, if statement fall-throughs

Exit points: Return statements, Goto's, if statements*

Entry and exit points are not included in the basic block, with the exception of if statements because the conditional in an if statement has the opportunity for optimization.

8.2 Methods

The following methods perform optimizations on the Three-Address-Code lines generated by the IR. These Three-Address-Code lines are then converted to instructions that are in the following format,

LHS = RHS

LHS = R1 OP R2

RHS = R1 OP R2

Each optimization method uses these instructions to return an optimized RHS. The extent of each optimization's utility is described below:

- Constant Folding: Method that evaluates the RHS of an assignment.
For example,

RHS = R1 OP R2

w = 3 + 5

After one pass,

$w = 8$

- Constant Propagation: Method that substitutes the value of known variables. For example,

$x = 5$
 $y = 9 - x$

After one pass,

$x = 5$
 $y = 9 - 5$

- Algebraic Simplification: This method simplifies the "algebra" in the RHS of the 3-Address code. For example:

- If $OP = *$ and $R1$ or $R2 = 0$
 $x = 0 * 5 \quad x = 5 * 0$
 $x = 0 \quad x = 0$
- If $OP = -$ and $R1 = R2$
 $x = 7 - 7$
 $x = 0$
- If $OP = /$ and $R1 = R2$
 $x = 3 / 3$
 $x = 1$

- Identity Removal: This method also simplifies the RHS of the expression but where there are variables. For example:

- If $OP = +$ or $-$ and $R1$ or $R2 = 0$
 $x = y - 0 \quad x = 0 + y$
 $x = y \quad x = y$
- If $OP = *$ or $/$ and $R1$ or $R2 = 1$
 $x = y / 1 \quad x = 1 * y$
 $x = y \quad x = y$
- If $OP = /$ and $R1 = R2$
 $x = 3 / 3$
 $x = 1$

- Dead Code Removal: Removes lines that assign to a variable that is never used in the RHS of any of the 3-Address code lines.

Unfortunately, the current dead code removal scheme cannot be used in the current version due to complications in the optimization design. In the current design, return statements are not included in a basic block, so a variable that is used only in a return statement could be removed by this optimization, which is incorrect.

It should be noted that, while these optimizations do apply to if statement conditionals, the conditionals themselves are limited. For example, if " $i < 10$ " were optimized to " $0 < 10$ ", the conditional would be left in this state.

9 Assembly Generation

9.1 Calling Convention

The compiler generates assembly in x86 AT&T and uses the cdecl calling convention. When a function is called, the function arguments are placed on the stack in reverse order, which allows the callee to reference the input arguments from the offset of EBP. When the function assembly is being generated, the compiler will count the number of local variables that the function uses, and will insert the correct size to decrement the stack pointer before any instructions are generated. Local variables are also referenced in relation to the EBP. On function exit, *leave*

9.2 Register Allocation

In terms of register allocation, the approach taken was a very simple one. During assembly generation, there is a list of registers that are currently not being used. When one is needed, a function is called to return the first free register, which is then marked as being in use. When the register is done being used in the context that it was called, the register is freed for later use.