

A Grammar for the C- Programming Language (Version S21)

January 19, 2021

1 Introduction

This is a grammar for the Spring 2021 semester's C- programming language. This language is very similar to C and has a lot of features in common with a real-world programming language. There are also some real differences between C and C-. For instance the declaration of procedure arguments, the loops that are available, what constitutes the body of a procedure, assignment is NOT a simple expression, array operators, etc. Also because of time limitations this language unfortunately does not have any heap related structures. It would be great to do a lot more, but we'll save for a second semester of compilers ☺. NOTE: this grammar is not a Bison grammar! You'll have to fix that as part of your assignment.

For the grammar that follows here are the types of the various elements by type font or symbol:

- **Keywords are in this type font.**
- **TOKEN CLASSES ARE IN THIS TYPE FONT.**
- *Nonterminals are in this type font.*
- The symbol ϵ means the empty string in a CS grammar sense.
- The token `|` is denoted `'|'` to distinguish it from the alternation symbol in the defining grammar.

1.1 Some Token Definitions

- letter = a | ... | z | A | ... | Z
- digit = 0 | ... | 9
- letdig = digit | letter
- **ID** = letter letdig*
- **NUMCONST** = digit⁺
- **CHARCONST** = is a representation for a single character by placing that character in **single quotes**. A backslash is an escape character. Any character preceded by a backslash is interpreted as that character. For example `\x` is the letter x, `\'` is a single quote, `\\` is a single backslash. There are **only two exceptions** to this rule: `\n` is a newline character and `\0` is the null character.
- **STRINGCONST** = any series of zero or more characters enclosed by **double quotes**. A backslash is an escape character. Any character preceded by a backslash is interpreted as that

character without meaning to the string syntax. For example `\x` is the letter x, `\"` is a double quote, `\'` is a single quote, `\\` is a single backslash. There are **only two exceptions** to this rule: `\n` is a newline character and `\0` is the null character. The string constant can be an empty string: a string of length 0. All string constants are terminated by the first unescaped double quote. String constants **must be entirely contained on a single line**, that is, they contain no unescaped newlines!

- **White space** (a sequence of blanks and tabs) is ignored. Whitespace may be required to separate some tokens in order to get the scanner not to collapse them into one token. For example: “intx” is a single **ID** while “int x” is the type **int** followed by the **ID** x. The scanner, by its nature, is a greedy matcher.
- **Comments** are ignored by the scanner. Comments begin with `//` and run to the end of the line.
- All **keywords** are in lowercase. You need not worry about being case independent since not all lex/flex programs make that easy.

2 The Grammar

1. $program \rightarrow declList$

2. $declList \rightarrow declList\ decl \mid decl$

3. $decl \rightarrow varDecl \mid funDecl$

4. $varDecl \rightarrow typeSpec\ varDeclList ;$

5. $scopedVarDecl \rightarrow \textbf{static}\ typeSpec\ varDeclList ; \mid typeSpec\ varDeclList ;$

6. $varDeclList \rightarrow varDeclList , varDeclInit \mid varDeclInit$

7. $varDeclInit \rightarrow varDeclId \mid varDeclId : simpleExp$

8. $varDeclId \rightarrow \textbf{ID} \mid \textbf{ID} [\textbf{NUMCONST}]$

9. $typeSpec \rightarrow \textbf{int} \mid \textbf{bool} \mid \textbf{char}$

10. $funDecl \rightarrow typeSpec\ \textbf{ID} (parms) stmt \mid \textbf{ID} (parms) stmt$

11. $parms \rightarrow parmList \mid \epsilon$

12. $parmList \rightarrow parmList ; parmTypeList \mid parmTypeList$

13. $parmTypeList \rightarrow typeSpec\ parmIdList$

14. $\text{parmIdList} \rightarrow \text{parmIdList}, \text{parmId} \mid \text{parmId}$

15. $\text{parmId} \rightarrow \mathbf{ID} \mid \mathbf{ID} []$

16. $\text{stmt} \rightarrow \text{expStmt} \mid \text{compoundStmt} \mid \text{selectStmt} \mid \text{iterStmt} \mid \text{returnStmt} \mid \text{breakStmt}$

17. $\text{expStmt} \rightarrow \text{exp}; \mid ;$

18. $\text{compoundStmt} \rightarrow \{ \text{localDecls stmtList} \}$

19. $\text{localDecls} \rightarrow \text{localDecls scopedVarDecl} \mid \epsilon$

20. $\text{stmtList} \rightarrow \text{stmtList stmt} \mid \epsilon$

21. $\text{selectStmt} \rightarrow \mathbf{if} \text{ simpleExp} \mathbf{then} \text{ stmt} \mid \mathbf{if} \text{ simpleExp} \mathbf{then} \text{ stmt} \mathbf{else} \text{ stmt}$

22. $\text{iterStmt} \rightarrow \mathbf{while} \text{ simpleExp} \mathbf{do} \text{ stmt} \mid \mathbf{for} \text{ ID} = \text{iterRange} \mathbf{do} \text{ stmt}$

23. $\text{iterRange} \rightarrow \text{simpleExp} \mid \text{simpleExp} \mathbf{to} \text{ simpleExp} \mid \text{simpleExp} \mathbf{to} \text{ simpleExp} \mathbf{by} \text{ simpleExp}$

24. $\text{returnStmt} \rightarrow \mathbf{return}; \mid \mathbf{return} \text{ exp};$

25. $\text{breakStmt} \rightarrow \mathbf{break};$

26. $\text{exp} \rightarrow \text{mutable} = \text{exp} \mid \text{mutable} += \text{exp} \mid \text{mutable} -= \text{exp} \mid \text{mutable} *= \text{exp} \mid \text{mutable} /= \text{exp} \mid \text{mutable} ++ \mid \text{mutable} -- \mid \text{simpleExp}$

27. $\text{simpleExp} \rightarrow \text{simpleExp} \mathbf{or} \text{ andExp} \mid \text{andExp}$

28. $\text{andExp} \rightarrow \text{andExp} \mathbf{and} \text{ unaryRelExp} \mid \text{unaryRelExp}$

29. $\text{unaryRelExp} \rightarrow \mathbf{not} \text{ unaryRelExp} \mid \text{relExp}$

30. $\text{relExp} \rightarrow \text{minmaxExp} \text{ relop} \text{ minmaxExp} \mid \text{minmaxExp}$

31. $\text{relop} \rightarrow <= \mid < \mid > \mid >= \mid == \mid !=$

32. $\text{minmaxExp} \rightarrow \text{minmaxExp} \text{ minmaxop} \text{ sumExp} \mid \text{sumExp}$

33. $\text{minmaxop} \rightarrow :>: \mid :<:$

34. $\text{sumExp} \rightarrow \text{sumExp} \text{ sumop} \text{ mulExp} \mid \text{mulExp}$

35. $\text{sumop} \rightarrow + \mid -$

36. $\text{mulExp} \rightarrow \text{mulExp} \text{ mulop} \text{ unaryExp} \mid \text{unaryExp}$

37. $\text{mulop} \rightarrow * \mid / \mid \%$

38. $unaryExp \rightarrow unaryop\ unaryExp \mid factor$
39. $unaryop \rightarrow - \mid * \mid ?$
40. $factor \rightarrow immutable \mid mutable$
41. $mutable \rightarrow \mathbf{ID} \mid \mathbf{ID} [exp]$
42. $immutable \rightarrow (exp) \mid call \mid constant$
43. $call \rightarrow \mathbf{ID} (args)$
44. $args \rightarrow argList \mid \epsilon$
45. $argList \rightarrow argList , exp \mid exp$
46. $constant \rightarrow \mathbf{NUMCONST} \mid \mathbf{CHARCONST} \mid \mathbf{STRINGCONST} \mid \mathbf{true} \mid \mathbf{false}$

3 Semantic Notes

- The only numbers are **ints**.
- There is no conversion or coercion between types such as between **ints** and **bools** or **bools** and **ints**.
- There can only be one function with a given name. There is no function overloading. The function name space is the same as the variable name space so a function and a variable cannot have the same name in the same scope.
- There are min and max operators denoted **:<:** and **:>:** respectively.
- The unary asterisk takes an array as an argument and returns the size of the array.
- The **STRINGCONST** token translates to a fixed size **char** array.
- The logical operators **and** and **or** are NOT short cutting. Although it is easy to do, we have plenty of other stuff to implement.
- In if statements the **else** is associated with the most recent **if**. The above grammar allows for ambiguous associations between **else** and **if**. In your assignment you will have to fix this.
- Exps are evaluated in order consistent with operator associativity and precedence found in mathematics. Also, no reordering of operands is allowed.
- A char occupies the same space as an integer or bool. This is an artifact of the virtual machine.
- A string constant is a constant char array.

- Initialization of variables can only be done with expressions that are constant, that is, they are able to be evaluated to a constant at compile time. For this class, it is not necessary that you actually evaluate the constant expression at compile time. But you will have to keep track of whether the expression is constant. Type of variable and expression must match (see exception for char arrays below).
- Array assignment works. The source array is copied to the target array. If the target array is smaller the source array is trimmed. If the target array is larger all the source elements are copied and the remainder of the target is untouched. There is hardware support for this.
- To be clear, assignment of a string (char array) to a char array works as if it is any other array assignment. It will not overrun the end of the lhs array. If it is too short the remainder of the array is untouched.
- Passing of arrays is done by reference implemented as pointers. Functions cannot return an array, but they can modify the content of an array passed in.
- Array comparison works. There is hardware support for this.
- Assignments in expressions happen at the time the assignment operator is encountered in the order of evaluation. The value returned is value of the rhs of the assignment. Assignments include the `++` and `--` operator. That is, the `++` and `--` operator do NOT behave as in C or C++. NOTE: assignment does NOT occur in a *simpleExp* without enclosing parens.
- The initializing a char array to a string behaves like an array assignment to the whole array.
- Initializing an array to a scalar constant fills the array with that constant.
- Function return type is specified in the function declaration, however, if no type is given to the function in the declaration then it is assumed the function does not return a value. To aid discussion of this case, the type of the return value is said to be void, even though there is no **void** keyword for the type specifier.
- Code that exits a procedure without a **return** returns a 0 for a function returning **int** and **false** for a function returning **bool** and a blank for a function returning **char**.
- All variables, functions must be declared before use.
- `?n` generates a uniform random integer in the range 0 to $|n| - 1$ with the sign of n attached to the result. `?5` is a random number in the range 0 to 4. `?-5` is a random number in the range 0 to -4. `?0` is undefined. There is hardware support for this.
- There are two types of loops. The **while** loops on a boolean condition and the **for** loops on an integer index variable. The **for** itself introduces a new scope. The **for** takes an **ID** which will be declared as an integer variable in the scope of the **for**. The range of values of the variable are defined in increasing detail by the other optional keyword combinations. The combinations are an initial value, a range of values, or and range of values and step. Step may be positive or negative. For simplicity, if step size is not given it is always assumed to be +1.

4 An Example of C- Code

```
char zev[10]:"corgis";
char yurt[20];
int x:42, y:666;

int ant(int bat, cat[]; bool dog, elk; int fox; char gnu)
{
    int goat, hog[100];

    gnu = 'W';
    goat = hog[2] = 3**cat;    // hog is 3 times the size of array passed to cat
    if dog and elk or bat > cat[3] then dog = not dog;
    else fox++;
    if bat <= fox then {
        while dog do {
            static int hog;          // hog in new scope

            hog = fox;
            dog = fred(fox++, cat)>666;
            if hog>bat then break;
            else if fox!=0 then fox += 7;
        }
    }

    for i = 0 to *zev-1 do outputc(zev[i]); outnl();
    if zev > "dog" then outputs("bark");
    yurt = zev;
    yurt[3] = zev[*zev];

    return (fox+bat*cat[bat])/~fox;
}

// note that functions are defined using a statement
int max(int a, b) if a>b then return a; else return b;

// use the max operator
int max3(int a, b, c) return a :>: b :>: c;
```

Table 1: A table of all operators in the language to help with type checking and showing what operations are allowed. Note that C- supports = for all types of arrays. Initialization is listed as an operator for type checking purposes and could have been implemented as an operator but I chose not to for no good reason. But the type checking needs to work.

| Binary Operator | Operands | Return Type |
|-----------------|----------------------|-------------|
| initialization | equal types + arrays | N/A |
| and | bool,bool | bool |
| — | bool,bool | bool |
| == | equal types + arrays | bool |
| != | equal types + arrays | bool |
| <= | equal types + arrays | bool |
| < | equal types + arrays | bool |
| >= | equal types + arrays | bool |
| > | equal types + arrays | bool |
| = | equal types + arrays | type of lhs |
| += | int,int | int |
| -= | int,int | int |
| *= | int,int | int |
| /= | int,int | int |
| :>: | int,int | int |
| :<: | int,int | int |
| * | int,int | int |
| + | int,int | int |
| — | int,int | int |
| / | int,int | int |
| % | int,int | int |
| [] | array,int | type of lhs |
| Unary Operator | Operands | Return Type |
| -- | int | int |
| ++ | int | int |
| not | bool | bool |
| * | array | int |
| — | int | int |
| ? | int | int |

Table 2: A table of just the array operators.

| Binary Operator | Operands | Return Type |
|-----------------|------------------|-------------|
| == | equal types | bool |
| != | equal types | bool |
| <= | equal (int/char) | bool |
| < | equal (int/char) | bool |
| >= | equal (int/char) | bool |
| > | equal (int/char) | bool |
| = | equal types | type of lhs |
| [] | array,int | type of lhs |
| Unary Operator | Operands | Return Type |
| * | any | int |