

A Grammar for the C- Programming Language (Version S20)

January 21, 2020

1 Introduction

This is a grammar for the Spring 2020 semester's C- programming language. This language is very similar to C and has a lot of features in common with a real-world structured 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 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.

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.

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. $program \rightarrow declarationList$
2. $declarationList \rightarrow declarationList\ declaration \mid declaration$
3. $declaration \rightarrow varDeclaration \mid funDeclaration$

-
4. $varDeclaration \rightarrow typeSpecifier\ varDeclList ;$
 5. $scopedVarDeclaration \rightarrow scopedTypeSpecifier\ varDeclList ;$
 6. $varDeclList \rightarrow varDeclList , varDeclInitialize \mid varDeclInitialize$
 7. $varDeclInitialize \rightarrow varDeclId \mid varDeclId : simpleExpression$
 8. $varDeclId \rightarrow ID \mid ID [NUMCONST]$
 9. $scopedTypeSpecifier \rightarrow static\ typeSpecifier \mid typeSpecifier$
 10. $typeSpecifier \rightarrow int \mid bool \mid char$

-
11. $funDeclaration \rightarrow typeSpecifier\ ID (params) statement \mid ID (params) statement$
 12. $params \rightarrow paramList \mid \epsilon$
 13. $paramList \rightarrow paramList ; paramTypeList \mid paramTypeList$
 14. $paramTypeList \rightarrow typeSpecifier\ paramIdList$

15. $paramIdList \rightarrow paramIdList , paramId \mid paramId$

16. $paramId \rightarrow \mathbf{ID} \mid \mathbf{ID} []$

17. $statement \rightarrow expressionStmt \mid compoundStmt \mid selectionStmt \mid iterationStmt \mid returnStmt \mid breakStmt$

18. $expressionStmt \rightarrow expression ; \mid ;$

19. $compoundStmt \rightarrow \{ localDeclarations statementList \}$

20. $localDeclarations \rightarrow localDeclarations scopedVarDeclaration \mid \epsilon$

21. $statementList \rightarrow statementList statement \mid \epsilon$

22. $elsifList \rightarrow elsifList \mathbf{elsif} simpleExpression \mathbf{then} statement \mid \epsilon$

23. $selectionStmt \rightarrow \mathbf{if} simpleExpression \mathbf{then} statement elsifList \mid \mathbf{if} simpleExpression \mathbf{then} statement elsifList \mathbf{else} statement$

24. $iterationRange \rightarrow \mathbf{ID} = simpleExpression .. simpleExpression \mid \mathbf{ID} = simpleExpression .. simpleExpression : simpleExpression$

25. $iterationStmt \rightarrow \mathbf{while} simpleExpression \mathbf{do} statement \mid \mathbf{loop} \mathbf{forever} statement \mid \mathbf{loop} iterationRange \mathbf{do} statement$

26. $returnStmt \rightarrow \mathbf{return} ; \mid \mathbf{return} expression ;$

27. $breakStmt \rightarrow \mathbf{break} ;$

28. $expression \rightarrow mutable = expression \mid mutable += expression \mid mutable -= expression \mid mutable *= expression \mid mutable /= expression \mid mutable ++ \mid mutable -- \mid simpleExpression$

29. $simpleExpression \rightarrow simpleExpression \mathbf{or} andExpression \mid andExpression$

30. $andExpression \rightarrow andExpression \mathbf{and} unaryRelExpression \mid unaryRelExpression$

31. $unaryRelExpression \rightarrow \mathbf{not} unaryRelExpression \mid relExpression$

32. $relExpression \rightarrow sumExpression relOp sumExpression \mid sumExpression$

33. $relOp \rightarrow <= \mid < \mid > \mid >= \mid == \mid !=$

34. $sumExpression \rightarrow sumExpression sumOp mulExpression \mid mulExpression$

35. $sumOp \rightarrow + \mid -$

36. $mulExpression \rightarrow mulExpression\ mulop\ unaryExpression \mid unaryExpression$
37. $mulop \rightarrow * \mid / \mid \%$
38. $unaryExpression \rightarrow unaryop\ unaryExpression \mid factor$
39. $unaryop \rightarrow - \mid * \mid ?$
40. $factor \rightarrow immutable \mid mutable$
41. $mutable \rightarrow \mathbf{ID} \mid mutable\ [expression]$
42. $immutable \rightarrow (expression) \mid call \mid constant$
43. $call \rightarrow \mathbf{ID}\ (args)$
44. $args \rightarrow argList \mid \epsilon$
45. $argList \rightarrow argList, expression \mid expression$
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 unary asterisk is the only unary operator that takes an array as an argument. It takes an array 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**.
- **elsif** is treated as if it is an **else** containing the **if** test and all the immediately following **elsif**'s. The rule of matching the **else** is associated with the most recent **if** applies here as a result.
- **loop** with a range creates a new scope with the **ID** declared as a variable in that scope. The from, to, and by values for range are computed once before the loop begins and stored in non-visible variables related to the loop and stored in the scope of the loop.

- Expressions 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.
- A string is a constant char array.
- Initialization of variables can only be 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 only the elements in the target corresponding to the source elements change. Array comparison doesn't work natively. We just don't have time. Passing of arrays is done by reference. Functions cannot return an array, but they can modify the content of an array passed in.
- 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++.
- Assignment of a string (char array) to a char array. This simply assigns all of the chars in the rhs array into the lhs array. It will not overrun the end of the lhs array. If it is too short it will pad the lhs array with null characters which are equivalent to zeroes.
- The initializing a char array to a string behaves like an array assignment to the whole array. ~~The second initializing case for a char array is to initialize it to a char (not a char array). This will fill the array with copies of the given character. By the way, this is an illegal assignment.~~
- 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 an 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. $?x$ for array x gives a random element from the array x .

4 An Example of C- Code

```
char zev[10]:"corgis";
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 = !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;
        }
    }

    loop i=1..10:3 do { // i is an int local to the loop
        if x==1 then cat[i] = bat;
        elsif (x==2) then cat[i] = bat%17;
        elsif (x==3) then cat[i] = 78;
        else x++;
    }

    loop forever if x>333 then break; else x++;

    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;
```

Table 1: A table of all operators in the language. Note that C- supports = for all types of arrays. It does not support relative testing: \geq , \leq , $>$, $<$ for any arrays.

Operator	Arguments	Return Type
initialization	equal,string	N/A
initialization	equal	N/A
not	bool	bool
and	bool,bool	bool
or	bool,bool	bool
==	equal types	bool
!=	equal types	bool
<=	int,int	bool
<	int,int	bool
>=	int,int	bool
>	int,int	bool
<=	char,char	bool
<	char,char	bool
>=	char,char	bool
>	char,char	bool
=	equal types incl. arrays	type of lhs
+=	int,int	int
-=	int,int	int
*=	int,int	int
/=	int,int	int
--	int	int
++	int	int
*	any array	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