

Computer Engineering 175

Phase II: Syntax Analysis

“Grammar, which knows how to control even kings.”
Molière, *Les Femmes Savantes*

1 Overview

In this assignment, you will write a recursive-descent parser for the Simple C language. This assignment is worth 20% of your project grade. Your program is due at 11:59 pm, Sunday, January 29th.

2 Syntactic Structure

The following rules constitute the syntax rules for Simple C:

<i>translation-unit</i>	→	ϵ <i>global-declaration translation-unit</i> <i>function-definition translation-unit</i>
<i>global-declaration</i>	→	<i>specifier global-declarator-list ;</i>
<i>global-declarator-list</i>	→	<i>global-declarator</i> <i>global-declarator , global-declarator-list</i>
<i>global-declarator</i>	→	<i>pointers id</i> <i>pointers id ()</i> <i>pointers id [num]</i>
<i>pointers</i>	→	ϵ <i>* pointers</i>
<i>specifier</i>	→	int char long void
<i>function-definition</i>	→	<i>specifier pointers id (parameters) { declarations statements }</i>
<i>parameters</i>	→	void <i>parameter-list</i>
<i>parameter-list</i>	→	<i>parameter</i> <i>parameter , parameter-list</i>
<i>parameter</i>	→	<i>specifier pointers id</i>
<i>declarations</i>	→	ϵ <i>declaration declarations</i>
<i>declaration</i>	→	<i>specifier declarator-list ;</i>
<i>declarator-list</i>	→	<i>declarator</i> <i>declarator , declarator-list</i>

<i>declarator</i>	→	<i>pointers id</i> <i>pointers id [num]</i>
<i>statements</i>	→	ϵ <i>statement statements</i>
<i>statement</i>	→	{ <i>declarations statements</i> } return <i>expression</i> ; while (<i>expression</i>) <i>statement</i> for (<i>assignment</i> ; <i>expression</i> ; <i>assignment</i>) <i>statement</i> if (<i>expression</i>) <i>statement</i> if (<i>expression</i>) <i>statement</i> else <i>statement</i> <i>assignment</i> ;
<i>assignment</i>	→	<i>expression</i> = <i>expression</i> <i>expression</i>
<i>expression</i>	→	<i>expression</i> <i>expression</i> <i>expression</i> && <i>expression</i> <i>expression</i> == <i>expression</i> <i>expression</i> != <i>expression</i> <i>expression</i> <= <i>expression</i> <i>expression</i> >= <i>expression</i> <i>expression</i> < <i>expression</i> <i>expression</i> > <i>expression</i> <i>expression</i> + <i>expression</i> <i>expression</i> - <i>expression</i> <i>expression</i> * <i>expression</i> <i>expression</i> / <i>expression</i> <i>expression</i> % <i>expression</i> ! <i>expression</i> - <i>expression</i> & <i>expression</i> * <i>expression</i> sizeof <i>expression</i> <i>expression</i> [<i>expression</i>] id (<i>expression-list</i>) id () id num string character (<i>expression</i>)
<i>expression-list</i>	→	<i>expression</i> <i>expression</i> , <i>expression-list</i>

3 Assignment

You will write a parser for Simple C, using the given grammar as a starting point. Unfortunately, the given expression grammar is ambiguous. Therefore, you must first disambiguate the grammar without changing the language accepted. To help you in your task, Table 1 shows the precedence and associativity of operators in Simple C.

Operators	Associativity	Arity	Output
[]	left	binary	index
& * ! - sizeof	right	unary	addr deref not neg sizeof
* / %	left	binary	mul div rem
+ -	left	binary	add sub
< > <= >=	left	binary	ltn gtn leq geq
== !=	left	binary	eql neq
&&	left	binary	and
	left	binary	or

Table 1: Operator associativity and precedence.

To illustrate that your parser is working correctly, you will write the operator, as shown in Table 1, used in each expression to the **standard output** after you have matched that expression. For example, `a + b * c` would generate `mul` and then `add` because the multiplication is done before the addition. In contrast, `a + b - c` would generate `add` and then `sub` since addition and subtraction have the same precedence but are left associative.

Your program will only be given **syntactically correct** programs as input. However, it is strongly advised that you test your program against syntactically incorrect programs as a way of finding errors in your implementation.

4 Hints

First, you will need to modify your lexical analyzer to return separate tokens for each keyword and operator. The parser will call the lexer when it needs a token. For simplicity, use the ASCII character value of a single-character token (e.g., `'+'`, `'-'`, `'*'`), and create a `#define` or enumeration for multi-character tokens such identifiers, numbers, keywords, and operators (e.g., `ID`, `NUM`, `RETURN`, `AND`).

To implement a recursive-descent parser, you will need to eliminate left recursion and left-factor the given grammar. The first step involves the rules for expressions. You can simply extend the example given in the textbook, writing one function for each level of precedence. Start by writing a parser just for expressions (i.e., the start symbol would be *expression*) and test it on expressions.

Left-factoring needs to be performed at several obvious places (e.g., *declarator*) and two non-obvious places. At the global level, we cannot immediately tell if we have a function definition as opposed to a global declaration. This problem can be solved by left-factoring the grammar, combining the rules for *function-definition*, *global-declaration*, and *global-declarator-list*. In a parameter list, we cannot immediately tell if the keyword **void** designates an empty parameter list or is the opening specifier in a non-empty list. This problem can be solved by left-factoring the rules for *parameters*, *parameter-list*, and *parameter*.