Project 1: Lexical and Syntax Analyzer

CS215 Shanghai Jiao Tong University

October 12, 2014

1 Introduction

C and some C-liked languages are the dominant programming languages. In this project, you are required to design and implement a simplified compiler, including a lexical analyzer and a syntax analyzer, for a given programming language, namely SMALE, which is a simplified C-liked language containing only the core part of C language. You will learn how to incrementally design and implement the successive phases of the compilation processes using off-the-shelf generators. In this project, Linux environment (e.g., Ubuntu or CentOS) is required.

$\mathbf{2}$ Lexical Analyzer

In this section, you are going to write a lexical analyser. The lexical analyser reads the source codes of SMALE and separates them into tokens.

2.1Tokens

```
\Rightarrow /* integer <sup>1*</sup>/
\Rightarrow /* identifier <sup>2*</sup>/
INT
ID
SEMI
                          \Rightarrow :
COMMA
                          \Rightarrow ,
BINARYOP \Rightarrow /* binary operators<sup>3</sup> */
```

¹A sequence of digits or digits followed by "0x(0X)" or "0" without spaces. In addition, the value should be in the range of $(-2^{31}, 2^{31})$ A character string consisting of alphabetic characters, digits and the underscore. In

addition, digits can't be the first character.

 $^{^{3}}$ See section 2.2.

```
UNARYOP \Rightarrow /* unary operators<sup>4</sup> */
TYPE
                   \Rightarrow int
LP
                   \Rightarrow (
RP
                   \Rightarrow )
LB
                   \Rightarrow [
RB
LC
RC
{\tt STRUCT}
                   \Rightarrow struct
RETURN
                  \Rightarrow return
IF
                   \Rightarrow if
ELSE
                  \Rightarrow else
BREAK
                  \Rightarrow break
CONT
                   \Rightarrow continue
FOR
                   \Rightarrow for
```

2.2 Operators

Operators in SMALE are shown below.

		Precedence	Operator	Associativity	Description
		1	()	Left-to-right	Function call or parenthesis
					Array subscripting
DO	\leftarrow				Structure element selection by reference
		2	(-)	Right-to-left	Unary minus
			!		Logical NOT
	,		++		Prefix increment
. 1					Prefix decrement
∇			~		Bit NOT
	1	3	*	Left-to-right	Product
	1		/		Division
			%		Modulus
		4	+		Plus
			(-)		Binary minus
		5	<<		Shift left
			>>		Shift right
		6	>		Greater than
			>=		Not less than
			<		Less than
			<=		Not greater than

 $^{^4}$ See section 2.2.

7	==		Equal to
	! =		Not equal to
8	&		Bit AND
9	^		Bit XOR
10			Bit OR
11	&&		Logical AND
12			Logical OR
13	=	Right-to-left	Assign
	+=		+ and assign
	-=		- and assign
	*=		* and assign
	/=		/ and assign
	&=		& and assign
	^=		^ and assign
	=		and assign
	<<= >>=		<< and assign
	>>=		>> and assign

2.3 Flex

Flex is short for *fast lexical analyzer generator*, which is a free version of lex written in C. In this project, you can use flex to generate the lexical analyzer.

Here are some references about Flex:

- Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman, Compilers: Principles, Techniques, and Tools, Second Edition. Chapter 3.5.
- http://en.wikipedia.org/wiki/Flex_lexical_analyser.
- http://flex.sourceforge.net/manual/.

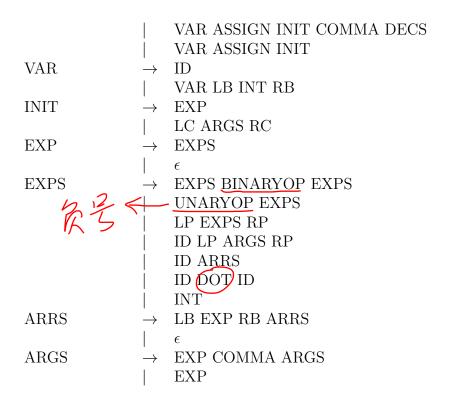
3 Syntax Analyzer

In this step, you are going to perform syntax analysis using Yacc.

3.1 Grammar⁵

⁵Input and output are not included, which will be provided in ??

PROGRAM \rightarrow EXTDEFS EXTDEFS EXTDEF EXTDEFS EXTDEF \rightarrow TYPE EXTVARS SEMI STSPEC SEXTVARS SEMI TYPE FUNC STMTBLOCK SEXTVARS IDID COMMA SEXTVARS ϵ **EXTVARS** VARVAR ASSIGN INIT VAR COMMA EXTVARS VAR ASSIGN INIT COMMA EXTVARS STRUCT ID LC SDEFS RC STSPEC STRUCT LC SDEFS RC STRUCT ID FUNC \rightarrow ID LP PARAS RP PARAS \rightarrow TYPE ID COMMA PARAS TYPE ID STMTBLOCK \rightarrow LC DEFS STMTS RC STMTS STMT STMTS STMT EXP SEMI STMTBLOCK RETURN EXP SEMI IF LP EXP RP STMT IF LP EXP RP STMT ELSE STMT FOR LP EXP SEMI EXP SEMI EXP RP STMT CONT SEMI BREAK SEMI TYPE DECS SEMI DEFS DEFS STSPEC SDECS SEMI DEFS TYPE SDECS SEMI SDEFS SDEFS SDECS ID COMMA SDECS ID DECS \rightarrow VAR VAR COMMA DECS



3.2 Notes and Hints

- The meanings of statements in SMALE is based on the meanings in C.
- The grammar given above may induce one or two reduce-reduce/shift-reduce conflicts, you need to assign precedence of some expressions manually to eliminate these conflicts. For example, "IF LP EXP RP STMT" should have lower precedence than "IF LP EXP RP STMT ELSE STMT".
- A number starts with '0x' or '0X' is a hexadecimal number, while a number starts with '0' is a octal number.
- Only integers and 1-dimensional array can be initialized.
- The dimension of arrays will be no more than 2.
- Struct can only contain int variables.
- The return type of a function can only be *int*.
- There are no "strange" statements such as a - b.

3.3 Yacc

Yacc is an LALR parser generator, which stands for *yet another compiler*-compiler. In this project, we can use yacc/bison (bison is another version of yacc) to generate a parser.

Here are some references about Yacc:

- Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman, Compilers: Principles, Techniques, and Tools, Second Edition. Chapter 4.9.
- http://en.wikipedia.org/wiki/Yacc.
- http://www.gnu.org/software/bison/manual/.

Submission Requirements

1. Pack the source files into a file named **StudentID-prj1.tar**. Please DO use "tar" command to compress your source files. Your source files should include:

File name	Description
smallc.l	Lex program
smallc.y	Yacc program
makefile	makefile
StudentID-report.pdf	project report

Table 4: File List.

- 2. Your analyzer will be tested by the following command:
 ./program "Source file name"
 Your analyzer is expected to read source codes from a source file, and output the lexical and syntax analysis result.
- 3. Please state clearly the purpose of each program at the head of the source codes, and comment your programs if necessary.
- 4. Send your StudentID-prj1.tar file to cs215.sjtu@gmail.com.
- 5. Due date: Midnight, Nov. 6, 2014.