

COP5621 Compiler Construction

Part I: Lexical Analyzer

Lexical Analyzer

In this assignment, you will write a lexical analyzer for the PASC language. The analyzer will consist of a scanner, written in Lex, and the routines to manage a lexical table.

Token Specification

In PASC, not all keywords are reserved words. The table in Page 3 defines the tokens that must be recognized, with their associated symbolic names. All multi-symbol tokens are separated by blanks, tabs, newlines, comments or delimiters.

Comments are enclosed in (* ... *), and cannot be nested. An identifier is a sequence of (upper or lower case) letters or digits, beginning with a letter. Upper and lower case are not distinguished (i.e. the identifier ABC is the same as Abc). There is no limit on the length of identifiers. However, you may impose limits on the total number of distinct identifiers and string lexemes and on the total number of characters in all distinct identifiers and string taken together. There should be no other limitation on the number of lexemes that the lexical analyzer will process.

An integer constant is an unsigned sequence of digits representing a 10-based number. A character constant is a character enclosed in single quotes (e.g., 'a' is a character constant a). A string constant is a sequence of characters surrounded by single quotes (e.g. 'Hello, world'). The internal representations of character constant and string constant are different in PASC. Hard-to-type or invisible characters can be represented in character and string constants by *escape sequences*; these sequences look like two characters, but represent only one. The escape sequences supported by the PASC language are \n for newline, \t for tab, \' for the single quote and \\ for the backslash itself. Any other character follows backslash is not treated as an escape sequence.

Token attributes

A unique identification of each token (integer aliased with the symbolic token name) must be returned by the lexical analyzer. In addition, the lexical analyzer must pass extra information about some token to the parser. This extra information is passed to the parser as a single value, namely an integer, through a global variable as described below. For integer constants, the numeric value of the constant is passed. For character constants, the numeric value of the character in the local set is passed. In order to allow other phases of the compiler to access the original identifier lexeme, the lexical analyzer passes an integer uniquely identifying an identifier (other than reserved words). String constants are treated in the same way, with a unique identifying number being passed. The unique identifying number, for both identifiers and string constants, should be an index (pointer) into a string table created by the lexical analyzer to record the lexemes.

Implementation

The central routine of the scanner is *yylex*, an integer function that returns a *token number*, indicating the type (identifier, integer constant, semicolon, etc.), of the next token in the input stream. In addition to the token type, *yylex* must set the global variables *yyline* and *yycolumn* to the line and column number at which that token appears, and – in the case of identifiers and constants –, put the extra information needed, as described above, into the global integer variable *yyval.semantic_value*. Lex will write *yylex* for you, using the patterns and rules defined in your lex input file (which should be called *lexer.l*). Your rules must include the code to maintain *yyline*, *yycolumn* and *yyval*.

Reserved words may be handled as regular expressions or stored as part of the id table. For example, reserve words may be pre-stored in the string table so your program can determine a reserve word from an identifier

by the section of the table in which the lexeme is found. Keywords that are not reserved have to be part of the id table. Efficiency should be a factor in the management of id table and the string table.

You are to write a routine `Lex_error` that takes a message and line and column numbers and reports an error, printing the message and indicating the position of the error. You need only print the line and column number to indicate the position.

The `#define` mechanism should be used to allow the lexical analyzer to return token numbers symbolically. In order to avoid using token names that are reserved or significant in C or in the parser, the token names have been specified for you in the subsequent table. The parser and the lexical analyzer must agree on the token numbers to ensure correct communication between them. The token number can be chosen by you, as the compiler writer, or, by default, by *Yacc*. Regardless of how token numbers are chosen, the end-marker must have token number 0 or negative, and thus, your lexical analyzer must return a 0 (or a negative) as token number upon reaching the end of input.

Error Handling

Your lexical analyzer should recover from all malformed lexemes, as well as such things as string constants that extended across a line boundary or comments that are never terminated.

Due date and materials to be handed in

The assignment is due Tuesday, Jan 29. You will do a demo at the due day. Test cases named `hello.pasc`, `sieve.pasc`, `hist.pasc`, `error.pasc` and `lexerror.pasc` can be found the class web page. A driver file, `driver.c`, that calls `yylex` and prints each token with its value as the input is scanned, can also be the class web page. Please use the driver to test your lex program.

Grading policy:

- Recognizing all tokens (in the right order) (60).
- Correct implementation of Id table and string table (10)
- Conversion of the special characters (5)
- Report errors with line and column numbers (10).
- non-recognizable characters error (5)
- non-closed comment error (5)
- string constant over the line boundary error (5)

Table 1: Token numbers to be returned by yylex

Token	Symbolic name	Token	Symbolic name
end of file	EOFnumber	else	ELSEnumber
;	SEMInumber	elsif	ELSIFnumber
:	COLONnumber	end	ENDnumber
,	COMMMANumber	endif	ENDIFnumber
.	DOTnumber	endloop	ENDLOOPnumber
(LPARENnumber	endrec	ENDRECnumber
)	RPARENnumber	exit	EXITnumber
<	LTnumber	for	FORnumber
>	GTnumber	forward	FORWARDnumber
=	EQnumber	function	FUNCTIONnumber
-	MINUSnumber	if	IFnumber
+	PLUSnumber	is	ISnumber
*	TIMESnumber	loop	LOOPnumber
..	DOTDOTnumber	not	NOTnumber
:=	COLEQnumber	of	OFnumber
<=	LEnumber	or	ORnumber
>=	GEnumber	procedure	PROCEDUREnumber
<>	NEnumber	program	PROGRAMnumber
<i>identifier</i>	IDnumber	record	RECORDnumber
<i>integer constant</i>	ICONSTnumber	repeat	REPEATnumber
<i>char constant</i>	CCONSTnumber	return	RETURNnumber
<i>string constant</i>	SCONSTnumber	then	THENnumber
and	ANDnumber	to	TOnumber
array	ARRAYnumber	type	TYPEnumber
begin	BEGINnumber	until	UNTILnumber
constant	CONSTnumber	var	VARnumber
div	DIVnumber	while	WHILEnumber
downto	DOWNTOnumber		