

Compiler Project Source Language

Kenneth Sundberg

October 17, 2013

Dedicated to Dr. Steven J. Allan

Abstract

CPSL (Compiler Project Source Language) is the programming language implemented in the Compiler Construction course at Utah State University (CS 5300). CPSL is a strongly typed, static scoped, block structured, high level language in the spirit of Pascal.

Acknowledgement

With permission, this work is largely derived from the CPSL language created for CS 5300 by Dr. Allan. The mistakes are mine.

Contents

Abstract	v
Acknowledgement	vii
1 Definitions	1
2 Lexical Structure	3
2.1 Introduction	3
2.2 Keywords	3
2.3 Identifiers	3
2.3.1 Predefined Identifiers	4
2.4 Operators and Delimiters	4
2.5 Constants	4
2.5.1 Integer Constants	4
2.5.2 Character Constants	4
2.5.3 String Constants	5
2.5.4 Representing Characters in Character and String Constants	5
2.5.5 Comments	5
2.5.6 Blanks, Tabs, Spaces, and New Lines	5
3 Syntactic Structure	7
3.1 Declarations	7
3.1.1 Constant Declarations	7
3.1.2 Procedure and Function Declarations	7
3.1.3 Type Declarations	8
3.1.4 Variable Declarations	8
3.2 CPSL Statements	8
3.3 Expressions	9

List of Tables

2.1	Keywords of CPSL	3
2.2	Predefined Identifiers in CPSL	4
2.3	Operators and Delimiters of CPSL	4

Chapter 1

Definitions

The character set of CPSL is ASCII. In the following *letter* denotes any upper- or lower-case letter, and *digit* denotes any of the ten decimal digits 0 through 9.

The syntax of CPSL is expressed in a EBNF-like format as follows:

$$\begin{aligned} \langle Nonterminal \rangle &\rightarrow \langle Nonterminal \rangle \textbf{ token } \langle token \text{ with state } \rangle \\ &| \langle Optional \rangle? \langle Repeated \rangle^* \\ &| \langle empty \rangle \end{aligned}$$

Within this format non-terminal symbols will be formatted in mixed case and enclosed in angle brackets. Terminals without state, such as keywords, will be formatted in bold. Terminals with state, such as identifiers, will be formatted in lower case and enclosed in angle brackets. Symbol groups will be indicated with parenthesis. Optional symbols or symbol groups will be indicated with a question mark. Alternatives for a symbol or symbol group will be indicated with a vertical bar. The terminal *empty* in a production indicates an empty production.

Chapter 2

Lexical Structure

2.1 Introduction

The following is a detailed description of the lexical structure of CPSL. Any lexeme not described herein is an error and should result in an appropriate diagnostic message.

2.2 Keywords

Keywords are strings with predefined meaning. They can not be redefined by the user. These keywords are either all capitals or all lower-case, mixed-case variations are not keywords. For example, BEGIN and begin are keywords while Begin is not.

array	begin	chr	const	do	downto
else	elseif	end	for	forward	function
if	of	ord	pred	procedure	read
record	repeat	return	stop	succ	then
to	type	until	var	while	write

Table 2.1: Keywords of CPSL

2.3 Identifiers

Identifiers in CPSL consist of a letter, followed by a sequence of zero or more letters, digits, or underscores. Upper- and lower-case letters are considered **distinct**. There is **no limit** on the length of identifiers in CPSL.

2.3.1 Predefined Identifiers

CPSL has a small set of predefined identifiers. Unlike keywords, the meaning of these identifiers may be altered by the user.

Identifier		Meaning
integer	INTEGER	Basic Integer type
char	CHAR	Basic Character type
boolean	BOOLEAN	Basic Boolean type
string	STRING	Basic String type
true	TRUE	Boolean constant
false	FALSE	Boolean constant

Table 2.2: Predefined Identifiers in CPSL

2.4 Operators and Delimiters

CPSL has a set of symbols used as operators and delimiters. These characters are not a valid part of any keyword or identifier.

+	-	*	/	&		~	=
<>	<	<=	>	>=	.	,	:
;	()	[]	:=	%	

Table 2.3: Operators and Delimiters of CPSL

2.5 Constants

2.5.1 Integer Constants

Integer Constants in CPSL are of three forms, each denoting an integer in a different base.

- A sequence of digits beginning with a 0 is interpreted as an octal number
- A 0x followed by a sequence of digits is interpreted as a hexadecimal number
- Any other sequence of digits is interpreted as a decimal value

2.5.2 Character Constants

A character constant represents a *single* character and is enclosed in a pair of single quotes. A character constant may not be blank, the lexeme consisting of a pair of single quotes is an error.

2.5.3 String Constants

A string constant represents a multi-character sequence and is enclosed in a pair of double quotes. String constants may not contain double quotes. A string constant may be empty.

2.5.4 Representing Characters in Character and String Constants

The newline character (ASCII 10) may not appear between the single or double quotes of a character or string constant. Any printable character (ASCII 32 to 126 inclusive) can be represented as itself with the exception of `\`. Also such a constant can contain a `\` followed by any printable character. Such a `\`-escaped sequence is interpreted as the character after the `\` with the following exceptions:

`\n` line feed

`\r` carriage return

`\b` backspace

`\t` tab

`\f` form feed

2.5.5 Comments

A comment in CPSL begins with a `$` and continues to the end of the line.

2.5.6 Blanks, Tabs, Spaces, and New Lines

Blanks, tabs, spaces, and new lines (white space) delimit other tokens but are otherwise ignored. This does not hold inside character and string constants where such characters are interpreted as themselves.

Chapter 3

Syntactic Structure

The overall structure of a CPSL program is as follows:

$$\langle Program \rangle \rightarrow \langle ConstantDecl \rangle? \langle TypeDecl \rangle? \langle VarDecl \rangle? \\ (\langle ProcedureDecl \rangle \mid \langle FunctionDecl \rangle)^* \langle Block \rangle .$$

3.1 Declarations

3.1.1 Constant Declarations

$$\langle ConstantDecl \rangle \rightarrow \text{const } (\langle ident \rangle = \langle ConstExpression \rangle ;)^+$$

3.1.2 Procedure and Function Declarations

Procedure and Function Declarations have the following structure:

$$\langle ProcedureDecl \rangle \rightarrow \text{procedure } \langle ident \rangle (\langle FormalParameters \rangle) ; \text{forward} ; \\ \mid \text{procedure } \langle ident \rangle (\langle FormalParameters \rangle) ; \langle Body \rangle ;$$
$$\langle FunctionDecl \rangle \rightarrow \text{function } \langle ident \rangle (\langle FormalParameters \rangle) : \langle Type \rangle ; \text{forward} ; \\ \mid \text{function } \langle ident \rangle (\langle FormalParameters \rangle) : \langle Type \rangle ; \langle Body \rangle ;$$
$$\langle FormalParameters \rangle \rightarrow \langle empty \rangle \\ \mid \text{var? } \langle IdentList \rangle : \langle Type \rangle (; \text{var? } \langle IdentList \rangle : \langle Type \rangle)^*$$
$$\langle Body \rangle \rightarrow \langle ConstantDecl \rangle? \langle TypeDecl \rangle? \langle VarDecl \rangle? \text{Block}$$
$$\langle Block \rangle \rightarrow \text{begin } \langle StatementSequence \rangle \text{end}$$

Notice that procedure and function definitions *cannot* be nested. In other words, we have only a global environment and a local environment. Procedures and functions can only be defined in the global environment. In the local environment we can only define constants, types, and variables. Notice also the

provision for accommodating **forward** references to procedures and functions. Notice also that both procedures and functions requires parentheses in the definition even if there are no parameters. This is also true for their invocations.

3.1.3 Type Declarations

In CPSL there are four predefined types (see figure 2.2), and two type constructors. The type constructors are array and record.

$$\begin{aligned} \langle TypeDecl \rangle &\rightarrow \mathbf{type} (\langle ident \rangle = \langle Type \rangle ;) + \\ \langle Type \rangle &\rightarrow \langle SimpleType \rangle \\ &\quad | \langle RecordType \rangle \\ &\quad | \langle ArrayType \rangle \\ \langle SimpleType \rangle &\rightarrow \langle ident \rangle \\ \langle RecordType \rangle &\rightarrow \mathbf{record} (\langle IdentList \rangle : \langle Type \rangle ;)^* \mathbf{end} \\ \langle ArrayType \rangle &\rightarrow \mathbf{array} [\langle ConstExpression \rangle : \langle ConstExpression \rangle] \mathbf{of} \langle Type \rangle \\ \langle IdentList \rangle &\rightarrow \langle ident \rangle (, \langle ident \rangle)^* \end{aligned}$$

3.1.4 Variable Declarations

$$\langle VarDecl \rangle \rightarrow \mathbf{var} (\langle IdentList \rangle : \langle Type \rangle ;) +$$

3.2 CPSL Statements

Statements in CPSL have the following syntax:

$$\begin{aligned} \langle StatementSequence \rangle &\rightarrow \langle Statement \rangle (; \langle Statement \rangle)^* \\ \langle Statement \rangle &\rightarrow \langle Assignment \rangle \\ &\quad | \langle IfStatement \rangle \\ &\quad | \langle WhileStatement \rangle \\ &\quad | \langle RepeatStatement \rangle \\ &\quad | \langle ForStatement \rangle \\ &\quad | \langle StopStatement \rangle \\ &\quad | \langle ReturnStatement \rangle \\ &\quad | \langle ReadStatement \rangle \\ &\quad | \langle WriteStatement \rangle \\ &\quad | \langle ProcedureCall \rangle \\ &\quad | \langle NullStatement \rangle \\ \langle Assignment \rangle &\rightarrow \langle LValue \rangle := \langle Expression \rangle \end{aligned}$$

$\langle \text{IfStatement} \rangle \rightarrow \text{if } \langle \text{Expression} \rangle \text{ then } \langle \text{StatementSequence} \rangle (\text{elseif } \langle \text{Expression} \rangle \text{ then } \langle \text{StatementSequence} \rangle)^* (\text{else } \langle \text{StatementSequence} \rangle)? \text{ end}$
 $\langle \text{WhileStatement} \rangle \rightarrow \text{while } \langle \text{Expression} \rangle \text{ do } \langle \text{StatementSequence} \rangle \text{ end}$
 $\langle \text{RepeatStatement} \rangle \rightarrow \text{repeat } \langle \text{StatementSequence} \rangle \text{ until } \langle \text{Expression} \rangle$
 $\langle \text{ForStatement} \rangle \rightarrow \text{for } \langle \text{ident} \rangle := \langle \text{Expression} \rangle (\text{to|downto} \langle \text{Expression} \rangle \text{ do } \langle \text{StatementSequence} \rangle \text{ end}$
 $\langle \text{StopStatement} \rangle \rightarrow \text{stop}$
 $\langle \text{ReturnStatement} \rangle \rightarrow \text{return } \langle \text{Expression} \rangle?$
 $\langle \text{ReadStatement} \rangle \rightarrow \text{read } (\langle \text{LValue} \rangle (, \langle \text{LValue} \rangle)^*)$
 $\langle \text{WriteStatement} \rangle \rightarrow \text{write } (\langle \text{Expression} \rangle (, \langle \text{Expression} \rangle)^*)$
 $\langle \text{ProcedureCall} \rangle \rightarrow \langle \text{ident} \rangle ((\langle \text{Expression} \rangle (, \langle \text{Expression} \rangle)^*)?)$
 $\langle \text{NullStatement} \rangle \rightarrow \langle \text{empty} \rangle$

3.3 Expressions

$\langle \text{Expression} \rangle \rightarrow \langle \text{Expression} \rangle \mid \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle \ \& \ \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle = \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle <> \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle <= \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle >= \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle < \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle > \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle + \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle - \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle * \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle / \langle \text{Expression} \rangle$
 $\mid \langle \text{Expression} \rangle \% \langle \text{Expression} \rangle$
 $\mid \sim \langle \text{Expression} \rangle$
 $\mid - \langle \text{Expression} \rangle$
 $\mid (\langle \text{Expression} \rangle)$
 $\mid \langle \text{ident} \rangle ((\langle \text{Expression} \rangle (, \langle \text{Expression} \rangle)^*)?)$
 $\mid \text{chr } (\langle \text{Expression} \rangle)$
 $\mid \text{ord } (\langle \text{Expression} \rangle)$
 $\mid \text{pred } (\langle \text{Expression} \rangle)$
 $\mid \text{succ } (\langle \text{Expression} \rangle)$
 $\mid \langle \text{LValue} \rangle$

$$\langle LValue \rangle \rightarrow \langle ident \rangle ((\cdot \langle ident \rangle) | ([\langle Expression \rangle]))^*$$

$$\begin{aligned} \langle ConstExpression \rangle \rightarrow & \langle ConstExpression \rangle | \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle \& \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle = \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle <> \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle <= \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle >= \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle < \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle > \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle + \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle - \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle * \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle / \langle ConstExpression \rangle \\ & | \langle ConstExpression \rangle \% \langle ConstExpression \rangle \\ & | \sim \langle ConstExpression \rangle \\ & | - \langle ConstExpression \rangle \\ & | (\langle ConstExpression \rangle) \\ & | \langle intconst \rangle \\ & | \langle charconst \rangle \\ & | \langle strconst \rangle \\ & | \langle ident \rangle \end{aligned}$$

To resolve ambiguities in this grammar the following suffices:

- Arithmetic and Boolean binary operators are left-associative
- Relational operators are non-associative
- Unary minus and Boolean not are right-associative
- Operators have the following precedence (decreasing order)

- Unary – (negation)
- * / %
- + –
- = <> < <= > >=
- ~
- &
- |