

TDT4205 Problem Set 2

Spring 2012

Answers are to be submitted to *itslearning*, by Wednesday, Feb. 15th, 20:00.

All submitted source code MUST be able to compile and run on asti. Grading is **pass/fail**. Please read the assignment guidelines on itslearning before starting to work on the assignment. Requests for clarifications can be posted on the itslearning forum.

What to turn in

When turning in assignments, please turn in two files:

- (your_username).answers.pdf : Answers to non-programming questions (Part 1)
- (your_username).code.zip,tar.gz,tgz : All your code for this assignment, including makefiles and other necessary code (Part 2)

PART 1 - Theory

Task 1.1

A subset of Pascal expressions (ref. <http://www.seas.gwu.edu/~hchoi/teaching/cs160d/pascal.pdf>) can be abstracted into the grammar

$$\begin{aligned} E &\rightarrow S \mid S \text{ r } S \\ S &\rightarrow T \mid I \text{ T } \mid S \text{ a } T \\ T &\rightarrow F \mid T \text{ m } F \\ F &\rightarrow i \mid n \\ I &\rightarrow p \mid m \end{aligned}$$

Tabulate the FIRST and FOLLOW sets for each nonterminal in this grammar in a table like Tab. 1.

Table 1: FIRST and FOLLOW sets for nonterminals

NT	FIRST(NT)	FOLLOW(NT)
E		
S		
T		
F		
I		

Task 1.2

1. What makes a bottom-up parser more powerful than a top-down parser?
2. What kind of parser does GNU Bison generate?
3. Why do we even need a compiler generator? Why not just hand-code the compilers, like other programs?

Task 1.3

1. What characterizes a shift-reduce parser? How is parsing done in practice?
2. Some grammars may result in conflicts in shift-reduce parsers. Describe the difference between shift-reduce conflicts and reduce-reduce conflicts. Give a small example of both conflicts.

PART 2 - Programming: Lexer and parser

The directory in the code archive `ps2_skeleton.tgz` begins a compiler for a slightly modified version of VSL (“Very Simple Language”), defined by Bennett[1].

The lexical structure of VSL is defined as follows:

- *Whitespace* consists of the characters `'\t'`, `'\n'` and `' '`. It is ignored after lexical analysis.
- *Comments* begin with the sequence `'//'`, and last until the next `'\n'` character. They are ignored after lexical analysis.
- *Reserved words* are `FUNC`, `PRINT`, `RETURN`, `CONTINUE`, `IF`, `THEN`, `ELSE`, `FI`, `WHILE`, `DO`, `DONE` and `VAR`.
- *Operators* are assignment (`:=`), and the basic arithmetic operators `'+'`, `'-'`, `'*'`, `'/'` and `'**'` (where the last one represents the power of-operator).
- *Numbers* are sequences of one or more decimal digits (`'0'` through `'9'`).
- *Strings* are sequences of arbitrary characters (except `'\n'`), enclosed in double quote characters (`'\"'`). It is an error to break a string across multiple lines.
- *Identifiers* are sequences of at least one letter followed by an arbitrary sequence of letters and digits. Letters are defined as the upper- and lower-case english alphabet (`'A'` through `'Z'` and `'a'` through `'z'`), as well as underscore (`'_'`). Digits are the decimal digits, as above.

The syntactic structure is given in the context-free grammar in Fig. 1.

Building the program combines `src/vslc.c`, `src/scanner.l` and `src/parser.y` into a binary called `bin/vslc`, which runs the scanner/parser pair on input from `stdin` (or optionally, a file). The scanner can dump a text representation of tokens and lexemes on `stderr` as it finds them, and the parser can produce a similar representation of syntax trees.

In the subdirectory `vsl_programs` is an example program, along with files of tokens/parse tree generated by a correct scanner/parser pair.

The structure in the `vslc` directory will be similar throughout subsequent problem sets, as the compiler takes shape. See the slide set from the 3rd recitation for an explanation of its construction, and notes on writing Lex/Yacc specifications.

$\text{program} \rightarrow \text{function_list}$
 $\text{function_list} \rightarrow \text{function} \mid \text{function_list } \text{function}$
 $\text{statement_list} \rightarrow \text{statement} \mid \text{statement_list } \text{statement}$
 $\text{print_list} \rightarrow \text{print_item} \mid \text{print_list } ', ' \text{ print_item}$
 $\text{expression_list} \rightarrow \text{expression} \mid \text{expression_list } ', ' \text{ expression}$
 $\text{variable_list} \rightarrow \text{variable} \mid \text{indexed_variable} \mid \text{variable_list } ', ' \text{ variable} \mid$
 $\text{variable_list } ', ' \text{ indexed_variable}$
 $\text{argument_list} \rightarrow \text{expression_list} \mid \epsilon$
 $\text{parameter_list} \rightarrow \text{variable_list} \mid \epsilon$
 $\text{declaration_list} \rightarrow \text{declaration_list } \text{declaration} \mid \epsilon$
 $\text{function} \rightarrow \text{FUNC } \text{variable} \text{ '(' } \text{parameter_list } \text{')' } \text{statement}$

$\text{statement} \rightarrow \text{assignment_statement} \mid \text{return_statement} \mid \text{print_statement} \mid \text{null_statement} \mid$
 $\text{if_statement} \mid \text{while_statement} \mid \text{block}$

$\text{block} \rightarrow \text{'{' } \text{declaration_list } \text{statement_list } \text{'}'}$
 $\text{assignment_statement} \rightarrow \text{variable } \text{ASSIGN } \text{expression} \mid$
 $\text{variable } \text{'[' } \text{expression } \text{']' } \text{ASSIGN } \text{expression}$
 $\text{return_statement} \rightarrow \text{RETURN } \text{expression}$
 $\text{print_statement} \rightarrow \text{PRINT } \text{print_list}$
 $\text{null_statement} \rightarrow \text{CONTINUE}$

$\text{if_statement} \rightarrow \text{IF } \text{expression } \text{THEN } \text{statement } \text{FI} \mid$
 $\text{IF } \text{expression } \text{THEN } \text{statement } \text{ELSE } \text{statement } \text{FI}$

$\text{while_statement} \rightarrow \text{WHILE } \text{expression } \text{DO } \text{statement } \text{DONE}$

$\text{expression} \rightarrow \text{expression } \text{'+' } \text{expression} \mid \text{expression } \text{'-' } \text{expression} \mid \text{expression } \text{'*'} \text{expression} \mid$
 $\text{expression } \text{'/' } \text{expression} \mid \text{'-' } \text{expression} \mid \text{expression } \text{POWER } \text{expression} \mid$
 $\text{'(' } \text{expression } \text{')' } \mid \text{integer} \mid \text{variable} \mid \text{variable } \text{'(' } \text{argument_list } \text{')' } \mid \text{variable } \text{'[' } \text{expression } \text{']'}$

$\text{declaration} \rightarrow \text{VAR } \text{variable_list}$
 $\text{variable} \rightarrow \text{IDENTIFIER}$
 $\text{indexed_variable} \rightarrow \text{variable } \text{'[' } \text{integer } \text{'}'}$
 $\text{integer} \rightarrow \text{NUMBER}$
 $\text{print_item} \rightarrow \text{expression} \mid \text{text}$
 $\text{text} \rightarrow \text{STRING}$

Figure 1: Context Free Grammar of VSL

Task 2.1

Complete the Lex scanner specification in `src/scanner.l` using the lexical specification, so that it properly tokenizes VSL programs.

Task 2.2

A `node_t` structure is defined in `src/tree.h`. Complete the auxiliary functions `node_init` and `node_finalize` in `src/tree.c`, so that they can initialize/free `node_t`-sized memory areas passed to them by their first argument. The function `destroy_subtree` should recursively remove the subtree below a given root node, while `node_finalize` should only remove the memory associated with a single node.

Task 2.3

Complete `src/parser.y` to include the VSL grammar, with semantic actions to construct the program's parse tree using the functions defined in Task 3. The top-level production should assign the root node to the globally accessible `node_t` pointer 'root' (declared in `src/parser.y`).

Task 2.4

The power of-operator (**) is supposed to be right-associative, and with a higher precedence than multiplication, but lower than unary minus. Set the associativity and operator precedence for ** in `parser.y`.

References

- [1] Bennett, J. P.: *Introduction to Compiling Techniques*, McGraw-Hill, 1990