

CPS2000 - Compiler Theory and Practice

Course Assignment (Part I) 2020/2021

Department of Computer Science, University of Malta
Sandro Spina

March 24, 2021

Instructions

- This is **an individual** assignment and carries **50%** of the final **CPS2000** grade.
- The submission deadline is **18th June 2021**. A soft-copy of the report and all related files (including code) must be uploaded to the VLE by midnight of the indicated deadline. Hard copies **are not** required to be handed in. Source and executable files must be archived into a single .zip file before uploading to the VLE. It is the student's responsibility to ensure that the uploaded .zip file is valid. The PDF report must be submitted separately through the Turnitin submission system on the VLE.
- A report describing how you designed and implemented the different tasks of the assignment **is required**. Tasks (1-5) for which no such information is provided in the report will **not be** assessed and therefore, no marks assigned.
- You are welcome to share ideas and suggestions. However, under no circumstances should code be shared among students. Please remember that plagiarism will not be tolerated; the final submission must be entirely your work.
- Your implementation can be carried out in either C++ or Java.
- You are to allocate approximately **40** hours for this assignment.

Description

In this assignment you are to develop a lexer, parser, semantic analyser and interpreter for the programming language - *TeaLang*. This part of the assignment is composed of three major components: i) design and implementation of a FSA-based table-driven lexer and hand-crafted top-down LL(k) parser, ii) visitor classes to perform XML generation, semantic analysis and interpreter execution on the abstract syntax tree (AST) produced by the top-down parser and iii) a report detailing how you designed and implemented the different tasks.

TeaLang is an expression-based strongly-typed programming language. The language has C-style comments, that is, `//...` for line comments and `/*...*/` for block comments. The language is case-sensitive and each function is expected to return a value. *TeaLang* has 4 primitive types: 'float', 'int', 'bool' and 'string'. Binary operators, such as '+', require that the operands have matching types; the language does not perform any implicit/automatic typecast. The following is a syntactically and semantically correct *TeaLang* program:

```
float Square(x:float) {
    return x*x;
}

bool XGreaterThanY(x:float , y:float) {
    let ans:bool = true;
    if (y > x) { ans = false; }
    return ans;
}

bool XGreaterThanYv2(x:float , y:float) {
    return x > y;
}

float AverageOfThree(x:float , y:float , z:float) {
    let total:float = x + y + z;
    return total / 3;
}

string JoinStr(s1:string , s2:string) {
    let s3:string = s1 + s2;
    return s3;
}

let x:float = 2.4;
let y:float = Square(2.5);
print y; //6.25
print XGreaterThanY(x, 2.3); //true
print XGreaterThanYv2(Square(1.5), y); //false
print AverageOfThree(x, y, 1.2); //3.28
print JoinStr(" Hello", " World"); //Hello World
```

The *TeaLang* programming language

The following rules describe the syntax of *TeaLang* in EBNF. Each rule has three parts: a left hand side (LHS), a right-hand side (RHS) and the ‘::=’ symbol separating these two sides. The LHS names the EBNF rule whereas the RHS provides a description of this name. Note that the RHS uses four control forms namely sequence, choice, option and repetition. In a sequence order is important and items appear left-to-right. The stroke symbol (... | ...) is used to denote choice between alternatives. One item is chosen from this list; order is not important. Optional items are enclosed in square brackets ([...]) indicating that the item can either be included or discarded. Repeatable items are enclosed in curly brackets ({ ... }); the items within can be repeated **zero** or more times. For example, a *Block* consists of zero or more *Statement* enclosed in curly brackets.

$\langle Letter \rangle$::= [A-Za-z]
$\langle Digit \rangle$::= [0-9]
$\langle Printable \rangle$::= [\x20-\x7E]
$\langle Type \rangle$::= ‘float’ ‘int’ ‘bool’ ‘string’
$\langle BooleanLiteral \rangle$::= ‘true’ ‘false’
$\langle IntegerLiteral \rangle$::= $\langle Digit \rangle$ { $\langle Digit \rangle$ }
$\langle FloatLiteral \rangle$::= $\langle Digit \rangle$ { $\langle Digit \rangle$ } ‘.’ $\langle Digit \rangle$ { $\langle Digit \rangle$ }
$\langle StringLiteral \rangle$::= “” { $\langle Printable \rangle$ } “”
$\langle Literal \rangle$::= $\langle BooleanLiteral \rangle$ $\langle IntegerLiteral \rangle$ $\langle FloatLiteral \rangle$ $\langle StringLiteral \rangle$
$\langle Identifier \rangle$::= (‘_’ $\langle Letter \rangle$) { ‘_’ $\langle Letter \rangle$ $\langle Digit \rangle$ }
$\langle MultiplicativeOp \rangle$::= ‘*’ ‘/’ ‘and’
$\langle AdditiveOp \rangle$::= ‘+’ ‘-’ ‘or’
$\langle RelationalOp \rangle$::= ‘<’ ‘>’ ‘==’ ‘!=’ ‘<=’ ‘>=’
$\langle ActualParams \rangle$::= $\langle Expression \rangle$ { ‘,’ $\langle Expression \rangle$ }
$\langle FunctionCall \rangle$::= $\langle Identifier \rangle$ ‘(’ [$\langle ActualParams \rangle$] ‘)’
$\langle SubExpression \rangle$::= ‘(’ $\langle Expression \rangle$ ‘)’
$\langle Unary \rangle$::= (‘-’ ‘not’) $\langle Expression \rangle$
$\langle Factor \rangle$::= $\langle Literal \rangle$ $\langle Identifier \rangle$ $\langle FunctionCall \rangle$ $\langle SubExpression \rangle$ $\langle Unary \rangle$

$\langle Term \rangle ::= \langle Factor \rangle \{ \langle MultiplicativeOp \rangle \langle Factor \rangle \}$
 $\langle SimpleExpression \rangle ::= \langle Term \rangle \{ \langle AdditiveOp \rangle \langle Term \rangle \}$
 $\langle Expression \rangle ::= \langle SimpleExpression \rangle \{ \langle RelationalOp \rangle \langle SimpleExpression \rangle \}$
 $\langle Assignment \rangle ::= \langle Identifier \rangle '=' \langle Expression \rangle$
 $\langle VariableDecl \rangle ::= 'let' \langle Identifier \rangle ':' \langle Type \rangle '=' \langle Expression \rangle$
 $\langle PrintStatement \rangle ::= 'print' \langle Expression \rangle$
 $\langle RtrnStatement \rangle ::= 'return' \langle Expression \rangle$
 $\langle IfStatement \rangle ::= 'if' '(' \langle Expression \rangle ')' \langle Block \rangle ['else' \langle Block \rangle]$
 $\langle ForStatement \rangle ::= 'for' '(' [\langle VariableDecl \rangle] ';' \langle Expression \rangle ';' [\langle Assignment \rangle] ')' \langle Block \rangle$
 $\langle WhileStatement \rangle ::= 'while' '(' \langle Expression \rangle ')' \langle Block \rangle$
 $\langle FormalParam \rangle ::= \langle Identifier \rangle ':' \langle Type \rangle$
 $\langle FormalParams \rangle ::= \langle FormalParam \rangle \{ ',' \langle FormalParam \rangle \}$
 $\langle FunctionDecl \rangle ::= \langle type \rangle \langle Identifier \rangle '(' [\langle FormalParams \rangle] ')' \langle Block \rangle$
 $\langle Statement \rangle ::= \langle VariableDecl \rangle ';' | \langle Assignment \rangle ';' | \langle PrintStatement \rangle ';' | \langle IfStatement \rangle | \langle ForStatement \rangle | \langle WhileStatement \rangle | \langle RtrnStatement \rangle ';' | \langle FunctionDecl \rangle | \langle Block \rangle$
 $\langle Block \rangle ::= '{' \{ \langle Statement \rangle \} '}'$
 $\langle Program \rangle ::= \{ \langle Statement \rangle \}$

Task Breakdown

Task 1 - Table-driven lexer

In this first task you are to develop the lexer for the *TeaLang* language. The lexer is to be implemented using the table-driven approach which simulates the DFA transition function of the *TeaLang* micro-syntax. The lexer should be able to report any lexical errors in the input program.

[Marks: 20%]

Task 2 - Hand-crafted LL(k) parser

In this task you are to develop a hand-crafted predictive parser for the *TeaLang* language. The Lexer and Parser classes interact through the function *GetNextToken()* which the parser uses to get the next valid token from the lexer. Note that for the vast majority of cases, the parser only needs to read one symbol of lookahead (k=1) in order to determine which production rule to use. The parser should be able to report any syntax errors in the input program. A successful parse of the input should produce an abstract syntax tree (AST) describing the structure of the program.

[Marks: 30%]

Task 3 - AST XML Generation Pass

In OOP programming, the Visitor design pattern is used to describe an operation to be performed on the elements of an object structure without changing the classes on which it operates. In our case this object structure is the AST (**not** the parse tree) produced by the parser in Task 2. For this task you are to implement a visitor class to output a properly indented XML representation of the generated AST. Please note that additional notes - with examples - on the Visitor design pattern are available on the VLE. Check those out before starting this task.

```
let x : float = 3.142 * 20.0765;

<Decl>
  <Var Type="float">x</Id>
  <BinExprNode Op="*>
    <FloatConst>3.142</FloatConst>
    <FloatConst>20.0765</FloatConst>
  </BinExprNode>
</Decl>
```

[Marks: 5%]

Task 4 - Semantic Analysis Pass

For this task, you are to implement another visitor class to traverse the AST and perform type-checking (e.g. checking that variables are assigned to appropriately typed expressions, variables are not declared multiple times in the same scope, etc.). In addition to the global program scope, local scopes are created whenever a block is entered and destroyed when control leaves the block. Note that blocks may be nested and that to carry out this task, it is essential to have a proper implementation of a symbol table. Your compiler should be able to report any semantic analysis errors resulting from the traversal on the AST. Note that whereas *TeaLang* allows for function definitions within local scope, for this task, function definitions need only be declared within the global scope. An important check carried out of the semantic analysis visitor is that of checking that a function always returns a value.

```
string OverUnder50(int : age) {
    if (age < 50) {
        return "Under Fifty";
    }
    else {
        return "Over Fifty";
    }
}

let x:int = 45;
while (x < 50) {
    print OverUnder50(x);  \\ "Under Fifty" x5
    x = x + 1;
}

print OverUnder50(x);    \\ "Over Fifty"
```

[Marks: 25%]

Task 5 - Interpreter Execution Pass

For this task, you are to implement another visitor class to traverse the AST and simulate an interpreter which executes the test program. The '*print*' <Expression> statement can be used in your test programs to output the value of <Expression> to the console and determine whether the computation carried out by the interpreter visitor is correct. Note that the symbol table now needs to be used to also store the values (in addition to type as was done with Task 4) of variables currently within scope.

[Marks: 20%]

Report

In addition to the source and class files, you are to write and submit a report. Remember that tasks 1 to 5 for which no information is provided in the report will not be assessed. In your report include any deviations from the original EBNF, the salient points on how you developed the lexer / parser / interpreter (and reasons behind any decisions you took) including semantic rules and code execution, and any sample *TeaLang* programs you developed for testing the outcome of your compiler. In your report, state what you are testing for, insert the program AST and the outcome of your test. As an example, the *TeaLang* source script below, computes the answer of a real number raised to an integer power:

```
//Function definition for Power
float Pow(x:float , n:int)
{
    let y : float = 1.0;           //Declare y and set it to 1.0
    if ( n>0 )
    {
        for (; n>0 ; n=n-1)
        {
            y = y * x;             //Assignment y = y * x;
        }
    }
    else
    {
        for (; n<0 ; n=n+1)
        {
            y = y / x;             //Assignment y = y / x;
        }
    }
    return y;                      //return y as the result
}

let x : float = Pow(2.1,10);

print x; //prints to console 1667.988
```