

TDT4165 - Programming Languages

Assignment 2: Introduction to Language Theory

Benjamin Zubača

September 19, 2024

Task 1: mdc

- a) In the function `Lex`, we utilize the built-in function `String.tokens` to extract lexemes from the input string, specifying whitespace as the delimiter for token separation.

```
1 fun {Lex Input}
2   {String.tokens Input & }
3 end
```

- b) → g) In the `Tokenize` function, we employ the `Map` function along with pattern matching to transform each lexeme into its corresponding record type based on its content. We first check each lexeme against predefined operators, and if none match, we attempt to convert the lexeme into a `number` record using `String.toInt` for the conversion.

```
1 fun {Tokenize Lexemes}
2   {Map Lexemes
3     fun {$ Lexeme}
4       case Lexeme
5       of "+" then operator(type:plus)
6       [] "-" then operator(type:minus)
7       [] "*" then operator(type:multiply)
8       [] "/" then operator(type:divide)
9       [] "p" then command(print)
10      [] "d" then command(duplicate)
11      [] "i" then command(negate)
12      [] "c" then command(clear)
13      else
14        try number({String.toInt Lexeme})
15        catch _ then raise "invalid lexeme" end
16      end
17    end
18  end}
19 end
```

- c) → g) In the `Interpret` function, we recursively process the list of tokens using a stack to manage operands and results, primarily through the `DoInterpret` function. Within this function, we evaluate specific commands and operators, adjusting the stack accordingly. When encountering an operator, we ensure that there are sufficient operands on the stack before performing the operation using the `Compute` function, which handles the arithmetic operations. Additional commands such as `print`, `duplicate`, `negate`, and `clear` are implemented using similar stack manipulation techniques. The `Unwrap` function simplifies the stack by extracting the values from each token, returning a list of the underlying numbers.

```

1 fun {Interpret Tokens}
2   {Unwrap {DoInterpret nil Tokens}}
3 end
4
5 fun {DoInterpret Stack Tokens}
6   case Tokens
7   of nil then Stack
8   [] Head|Tail then
9     case Head
10    of number(_) then {DoInterpret Head|Stack Tail}
11    [] operator(type:Op) then
12      if {Length Stack} < 2 then
13        raise "Not enough operands for operator" end
14      else
15        case Stack of Right|Left|Rest then
16          {DoInterpret {Compute Left Right Op}|Rest Tail}
17        end
18      end
19    [] command(print) then
20      {Show {Unwrap Stack}}
21      {DoInterpret Stack Tail}
22    [] command(duplicate) then
23      case Stack of Top|Rest then
24        {DoInterpret Top|Top|Rest Tail}
25      end
26    [] command(negate) then
27      case Stack of number(N)|Rest then
28        {DoInterpret number(~N)|Rest Tail}
29      end
30    [] command(clear) then {DoInterpret nil Tail}
31    end
32  end
33 end
34
35 fun {Compute Left Right Op}
36   case Op
37   of plus then number(Left.1 + Right.1)
38   [] minus then number(Left.1 - Right.1)
39   [] multiply then number(Left.1 * Right.1)
40   [] divide then number(Left.1 / Right.1)
41   else raise "Unknown operator" end
42   end
43 end
44
45 fun {Unwrap List}
46   {Map List fun {$ Token} Token.1 end}
47 end

```

A High-Level Description of the `mdc` Program

The `mdc` program starts by accepting a string input that represents a RPN expression. First, the `Lex` function is invoked to tokenize the input string, breaking it down into individual lexemes based on whitespace. This results in a list of tokens that represent either numbers or operators.

Next, the `Tokenize` function takes these lexemes and transforms them into structured records. Each token is categorized as either a numerical value (encapsulated in a `number` record) or an operator (encapsulated in an `operator` record).

During the interpretation phase, the `Interpret` function processes the list of token records recursively using a stack data structure. The core of this implementation lies within the `DoInterpret` function: As tokens are evaluated, numbers are pushed onto the stack, while operators initiate arithmetic operations. When an operator is encountered, the program checks for sufficient operands on the stack, executes the specified operation through the `Compute` function, and updates the stack with the result.

Additionally, the program includes commands for managing the stack, such as printing its current contents, duplicating the top element, negating the top element, and clearing the entire stack.

Once all tokens have been processed, the remaining stack contains the final computed values, which can be displayed or further manipulated as needed.

Task 2: Convert Postfix Notation Into an Expression Tree

- a) The function `ExpressionTreeInternal` processes tokens of a postfix expression recursively and builds the expression tree. It keeps track of numbers and intermediate results using an expression stack. When an operator is encountered, it pops the necessary operands from the stack, combines them into a tree node, and pushes the result back onto the stack.

```
1 fun {ExpressionTreeInternal Tokens ExpressionStack}
2   case Tokens
3   of nil then
4     case ExpressionStack
5     of Head|_ then Head
6     else raise "invalid expression" end
7   end
8   [] Head|Tail then
9     case Head
10    of number(N) then {ExpressionTreeInternal Tail N|ExpressionStack}
11    [] operator(type:Op) then
12      if {Length ExpressionStack} < 2 then
13        raise "not enough operands for operator" end
14      else
15        case ExpressionStack of Left|Right|Rest then
16          case Op
17          of plus then {ExpressionTreeInternal Tail plus(Left Right)|Rest}
18          [] minus then {ExpressionTreeInternal Tail minus(Left Right)|Rest}
19          [] multiply then {ExpressionTreeInternal Tail multiply(Left
20                           Right)|Rest}
21          [] divide then {ExpressionTreeInternal Tail divide(Left Right)|Rest}
22          else raise "unknown operator" end
23        end
24      end
25    else raise "invalid token" end
26  end
27 end
28 end
```

- b) The function `ExpressionTree` serves as a wrapper for `ExpressionTreeInternal`, initializing the function with an empty stack. It returns the root of the fully constructed expression tree once all tokens have been processed.

```
1 fun {ExpressionTree Tokens}
2   {ExpressionTreeInternal Tokens nil}
3 end
```

A High-Level Description of the Postfix to Expression Tree Conversion

The process of converting postfix notation into an expression tree follows a stack-based approach similar to evaluating postfix expressions. The algorithm processes each token in the input:

- Numbers are pushed onto the stack as leaf nodes.
- Operators pop two nodes from the stack (representing operands), create a new tree node where the operator is the parent and the two operands are its children, and then push the resulting subtree back onto the stack.

This continues until all tokens are processed, leaving a single node on the stack, which represents the root of the expression tree.

Task 3: Theory

a) The grammar of the lexemes from Task 1 can be formalized as $\Gamma = (V, S, R, v_s)$, where

$$V = \{\varepsilon, \lambda, \nu, \omega, \kappa\}$$

$$S = \{+, -, *, /, p, d, i, c, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, \sqcup\}$$

$$\begin{aligned} R = \{ & (\varepsilon, \lambda), (\varepsilon, \lambda \sqcup \varepsilon), \\ & (\lambda, \nu), (\lambda, \omega), (\lambda, \kappa), \\ & (\nu, 0), (\nu, 1), (\nu, 2), (\nu, 3), (\nu, 4), (\nu, 5), (\nu, 6), (\nu, 7), (\nu, 8), (\nu, 9), (\nu, \nu\nu), \\ & (\omega, +), (\omega, -), (\omega, *), (\omega, /), \\ & (\kappa, p), (\kappa, d), (\kappa, i), (\kappa, c)\} \end{aligned}$$

$$v_s = \varepsilon.$$

Here, $v_s = \varepsilon$ represents the start variable, λ represents a single lexeme, ν represents an integer, ω represents an operator, κ represents a command, and \sqcup represents a whitespace symbol.

This grammar captures the structure of the lexemes in Task 1. An expression can be a single lexeme or multiple lexemes separated by whitespace. Integers can be single digits or multiple digits (represented by the recursive production $\nu \rightarrow \nu\nu$). Operators are the four basic arithmetic operations. Commands are single-letter instructions for print, duplicate, negate, and clear. Whitespace (\sqcup) is explicitly included to separate lexemes in multi-lexeme expressions.

b) The grammar of the records returned by the **ExpressionTree** function in Task 2 can be formalized as the following:

$$\begin{aligned} \langle expr \rangle & ::= \text{integer} \\ & \quad | \langle oper \rangle \\ \langle oper \rangle & ::= \langle plus \rangle \\ & \quad | \langle minus \rangle \\ & \quad | \langle multiply \rangle \\ & \quad | \langle divide \rangle \\ \langle plus \rangle & ::= \text{plus}(\langle expr \rangle, \langle expr \rangle) \\ \langle minus \rangle & ::= \text{minus}(\langle expr \rangle, \langle expr \rangle) \\ \langle multiply \rangle & ::= \text{multiply}(\langle expr \rangle, \langle expr \rangle) \\ \langle divide \rangle & ::= \text{divide}(\langle expr \rangle, \langle expr \rangle) \end{aligned}$$

- c) The grammar in step a) can be classified as context-free because each production rule contains a single non-terminal on the left-hand side, and the right-hand side consists of a sequence of terminals and/or non-terminals. For instance, the recursive production $\nu \rightarrow \nu\nu$ enables the generation of multi-digits integers. In contrast, this grammar cannot be considered regular, as regular grammars restrict rules to those where the left-hand non-terminal is replaced by at most one terminal and possibly one-terminal. The recursion in $\nu \rightarrow \nu\nu$ violates this constraint, making the grammar non-regular.

Similarly, the grammar in part (b) is also context-free, as it describes tree-like structures for mathematical expressions. For example,

$\langle plus \rangle ::= \text{plus}(\langle expr \rangle, \langle expr \rangle).$

This rule follows the context-free pattern, with a single non-terminal on the left-hand side and a combination of terminals and/or non-terminals on the right-hand side. Like the grammar in part a), this one cannot be classified as regular either, because it involves non-linear structures and recursive patterns. Regular grammars, being more restrictive, only allow productions with limited right-hand sides, typically a single terminal followed by at most one non-terminal.