

**CS 316 E6TBA Queens College CUNY**  
**Exercise Set #2**

1. Consider the following EBNF grammar:

$\langle \text{id} \rangle \rightarrow \langle \text{letter} \rangle \{ \langle \text{letter} \rangle \mid \langle \text{digit} \rangle \}$

$\langle \text{extended id} \rangle \rightarrow \langle \text{id} \rangle \{ \text{"\_"} \langle \text{letters and digits} \rangle \}$

$\langle \text{letters and digits} \rangle \rightarrow \{ \langle \text{letter} \rangle \mid \langle \text{digit} \rangle \}^+$

Construct a DFA that accepts the strings in  $\langle \text{extended id} \rangle$ .

2. Given is the following BNF grammar:

$\langle \text{comp op} \rangle \rightarrow "<" \mid "<=" \mid ">" \mid ">=" \mid "=" \mid "!="$

$\langle \text{arith op} \rangle \rightarrow + \mid - \mid * \mid /$

$\langle \text{paren} \rangle \rightarrow "(" \mid ")"$

$\langle \text{token} \rangle \rightarrow \langle \text{comp op} \rangle \mid \langle \text{arith op} \rangle \mid \langle \text{paren} \rangle$

Construct a DFA that accepts the strings in  $\langle \text{token} \rangle$ .

3. Given is the following EBNF grammar:

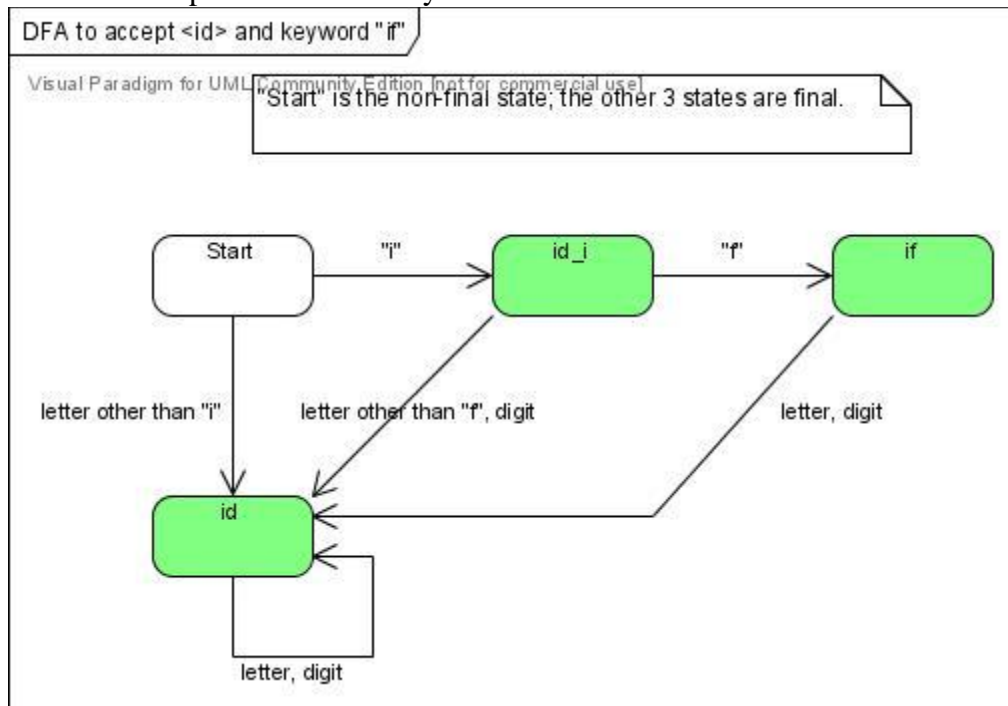
$\langle \text{float} \rangle \rightarrow [+ \mid -] \{ \langle \text{digit} \rangle \}^+ "." \{ \langle \text{digit} \rangle \}^+ [ (E \mid e) [+ \mid -] \{ \langle \text{digit} \rangle \}^+ ]$

Construct a DFA that accepts the strings in  $\langle \text{float} \rangle$ . Note that according to this grammar, the integer part may be empty.

4. Building keyword recognition into DFAs is not conceptually difficult but considerably increases the number of states and transitions. The key fact about usual keywords is that an identifier may contain an initial segment or the whole of a keyword; e.g., "whil", "whil5", "whilea", etc. Based on this observation the following method can be used:

- Create a linear sequence of transitions for each keyword; any non-empty, proper initial segment of each keyword is accepted by a final state as an identifier.
- Whenever the input letter or digit breaks a keyword, the transition branches off to the  $\langle \text{id} \rangle$  state.
- The final state for each keyword has the transition to the  $\langle \text{id} \rangle$  state on the letters and digits.

The following is a DFA to accept  $\langle \text{id} \rangle$  and the keyword "if":



Modify this DFA to accept additional keywords "else" and "while".

5. Consider the following EBNF grammar for `<int>` that defines a subset of the Java integer numerals.

`<digit>`  $\rightarrow 0 \mid 1 \mid \dots \mid 9$

`<nonZeroDigit>`  $\rightarrow 1 \mid \dots \mid 9$

`<hexDigit>`  $\rightarrow \text{<digit>} \mid a \mid b \mid \dots \mid f \mid A \mid B \mid \dots \mid F$

`<octalDigit>`  $\rightarrow 0 \mid 1 \mid \dots \mid 7$

`<int>`  $\rightarrow \text{<decimal>} \mid \text{<hexadecimal>} \mid \text{<octal>}$

`<decimal>`  $\rightarrow 0 \mid \text{<nonZeroDigit>} \{ \text{<digit>} \}$

`<hexadecimal>`  $\rightarrow ("0x" \mid "0X") \{ \text{<hexDigit>} \}^+$

`<octal>`  $\rightarrow 0 \{ \text{<octalDigit>} \}^+$

Construct a DFA that accepts the strings in `<int>`.

(Note: C++ integer numerals are almost identical; the difference is that "0" belongs to `<octal>`.)

6. Lexical analyzers are an example of pattern matching based on BNF grammars and finite automata, which has many other applications. Quite a number of modern script and programming languages provide built-in features or library functions to recognize patterns specified by EBNF grammars.

1. Suppose we want to find all file names containing an occurrence of "CS316" and ending with the extension ".txt". The "C" or "S" may be in lowercase. Give an EBNF grammar that defines this pattern.
2. Spam filters. Suppose we want to catch all email subject headers containing an occurrence of "pharmacy" where any of its characters may be in uppercase and any single ASCII character may be inserted between any of its characters. For example, "pHar%Ma?cy" and "p+hArMac!y". Give an EBNF Grammar that defines this pattern.

7. Consider the following grammar:

`<sequence>`  $\rightarrow "(" \text{<elements>} ")"$

`<elements>`  $\rightarrow \text{<element>} \{ "," \text{<element>} \}$

`<element>`  $\rightarrow \text{<id>} \mid \text{<sequence>}$

`<id>`  $\rightarrow \text{<letter>} \{ \text{<letter>} \mid \text{<digit>} \}$

Give a top-down parser in pseudo code for `<sequence>`, `<elements>`, and `<element>`. Presume that the function `getToken()` is given to extract the next token and assign it to the string variable *t*. The tokens are: `<id>`, `"("`, `)"`, and `","`.

8. Consider the following grammar for arithmetic expressions, with the unary `-` operator in `<primary>`.

`<E>`  $\rightarrow \text{<term>} \{ (+ \mid -) \text{<term>} \}$

`<term>`  $\rightarrow \text{<primary>} \{ ( * \mid / ) \text{<primary>} \}$

`<primary>`  $\rightarrow \text{<id>} \mid "(" \text{<E>} ")" \mid - \text{<primary>}$

Presume that the function `getToken()` is given to extract the next token and assign it to the string variable *t*. The tokens are: `<id>`, `"("`, `)"`, `+`, `-`, `*`, and `/`.

1. Give pseudo code for a top-down parser for `<E>`, `<term>`, and `<primary>`; include pseudo code to emit suitable stack-based instructions *push*, *add*, *sub*, *mul*, *div*, and *neg*.
2. Give the sequence of *push*, *add*, *sub*, *mul*, *div*, and *neg* instructions that would be generated from each of the following expressions:
  - a.  $X+Y+Z$
  - b.  $X-Y+Z-W$
  - c.  $X+Y*Z$
  - d.  $(X+Y)/(Z-W)$
  - e.  $X*-Y+Z/-W$

9. Consider the following grammar:

`<A>`  $\rightarrow (+ \mid -) \text{<B>} \text{<B>}$

`<B>`  $\rightarrow \text{<id>} \mid \text{<A>}$

Give pseudo code for a top-down parser for `<A>` and `<B>`. Presume that the function `getToken()` is given to extract the next token and assign it to the string variable *t*. The tokens are: `<id>`, `+`, and `-`.

**10.** In class we studied a top-down parser for the following grammar:

$\langle S \rangle \rightarrow \langle \text{assignment} \rangle \mid \{ \langle S \text{ List} \rangle \}$

$\langle \text{assignment} \rangle \rightarrow \langle \text{id} \rangle = \langle E \rangle ;$

$\langle S \text{ List} \rangle \rightarrow \{ \langle S \rangle \}^+$

Let us extend the production rule for  $\langle S \rangle$  with while-loops and conditionals:

$\langle S \rangle \rightarrow \langle \text{assignment} \rangle \mid \{ \langle S \text{ List} \rangle \} \mid \text{while } ( \langle B \rangle ) \langle S \rangle \mid \text{if } ( \langle B \rangle ) \langle S \rangle [ \text{else } \langle S \rangle ]$

Give pseudo code for a top-down parser for the extended  $\langle S \rangle$ ,  $\langle \text{assignment} \rangle$ , and  $\langle S \text{ List} \rangle$ . Presume that you are given the function  $B()$  to parse Boolean expressions  $\langle B \rangle$ ,  $E()$  to parse arithmetic expressions  $\langle E \rangle$ , and  $\text{getToken}()$  to extract the next token and assign it to the string variable  $t$ . The tokens are:  $\langle \text{id} \rangle$ , **if**, **else**, **while**,  $\{$ ,  $\}$ ,  $($ ,  $)$ ,  $=$ ,  $;$ . The top-down parser, if correctly constructed according to this grammar, should automatically build implicit parse trees in which each **else** matches the closest preceding unmatched **if**.