



Usman Institute of Technology

Department of Computer Science

Course Code: CS412

Course Title: Compiler Construction

Fall 2022

Lab 03

Objective:

This experiment introduces the students to implement the concept of regular expressions in programming. Also how to extract the regular expressions for a given DFA then do code.

Student Information

| | |
|--------------|--|
| Student Name | |
| Student ID | |
| Date | |

Assessment

| | |
|----------------|--|
| Marks Obtained | |
| Remarks | |
| Signature | |

Usman Institute of Technology
Department of Computer Science
CS412 – Compiler Construction

Lab 03

Instructions

- Come to the lab in time. Students who are late more than 10 minutes, will not be allowed to attend the lab.
 - Students have to perform the examples and exercises by themselves.
 - Raise your hand if you face any difficulty in understanding and solving the examples or exercises.
 - Lab work must be submitted on or before the submission date.
 - Do not copy the work of other students otherwise both will get zero marks.
-

1. Objective

This experiment introduces the students to implement the concept of regular expressions in programming. Also, how to extract the regular expressions for a given DFA then do code.

2. Labs Descriptions

What is Regular Expression?

Just as finite automata are used to recognize patterns of strings, regular expressions are used to generate patterns of strings. A regular expression is an algebraic formula whose value is a pattern consisting of a set of strings, called the language of the expression.

Operands in a regular expression can be:

characters from the alphabet over which the regular expression is defined.

variables whose values are any pattern defined by a regular expression.

epsilon which denotes the empty string containing no characters.

null which denotes the empty set of strings.

Operators used in regular expressions include:

Union: If R1 and R2 are regular expressions, then $R1 \mid R2$ (also written as $R1 \cup R2$ or $R1 + R2$) is also a regular expression.

$$L(R1 \mid R2) = L(R1) \cup L(R2).$$

Concatenation: If R1 and R2 are regular expressions, then $R1R2$ (also written as $R1.R2$) is also a regular expression.

$$L(R1R2) = L(R1) \text{ concatenated with } L(R2).$$

Kleene Closure: If R1 is a regular expression, then $R1^*$ (the Kleene closure of R1) is also a regular expression.

$$L(R1^*) = \epsilon \cup L(R1) \cup L(R1R1) \cup L(R1R1R1) \cup \dots$$

Lab 03: Design a code for regular expression to DFA and vice versa

Closure has the highest precedence, followed by concatenation, followed by union.

Examples

The set of strings over $\{0,1\}$ that end in 3 consecutive 1's.

$$(0 \mid 1)^* 111$$

The set of strings over $\{0,1\}$ that have at least one 1.

$$0^* 1 (0 \mid 1)^*$$

The set of strings over $\{0,1\}$ that have at most one 1.

$$0^* \mid 0^* 1 0^*$$

The set of strings over $\{A..Z,a..z\}$ that contain the word "main".

Let $\langle \text{letter} \rangle = A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$

$$\langle \text{letter} \rangle^* \text{main} \langle \text{letter} \rangle^*$$

The set of strings over $\{A..Z,a..z\}$ that contain 3 x's.

$$\langle \text{letter} \rangle^* x \langle \text{letter} \rangle^* x \langle \text{letter} \rangle^* x \langle \text{letter} \rangle^*$$

The set of identifiers in Pascal.

Let $\langle \text{letter} \rangle = A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$

Let $\langle \text{digit} \rangle = 0 \mid 1 \mid 2 \mid 3 \dots \mid 9$

$$\langle \text{letter} \rangle (\langle \text{letter} \rangle \mid \langle \text{digit} \rangle)^*$$

The set of real numbers in Pascal.

Let $\langle \text{digit} \rangle = 0 \mid 1 \mid 2 \mid 3 \dots \mid 9$

Let $\langle \text{exp} \rangle = 'E' \langle \text{sign} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle^* \mid \text{epsilon}$

Let $\langle \text{sign} \rangle = '+' \mid '-' \mid \text{epsilon}$

Let $\langle \text{decimal} \rangle = '.' \langle \text{digit} \rangle \langle \text{digit} \rangle^* \mid \text{epsilon}$

$$\langle \text{digit} \rangle \langle \text{digit} \rangle^* \langle \text{decimal} \rangle \langle \text{exp} \rangle$$

Unix Operator Extensions

Regular expressions are used frequently in Unix:

In the command line

Within text editors

In the context of pattern matching programs such as grep and egrep

To facilitate construction of regular expressions, Unix recognizes additional operators. These operators can be defined in terms of the operators given above; they represent a notational convenience only.

character classes: $[\langle \text{list of chars} \rangle]$

start of a line: '^'

end of a line: '\$'

wildcard matching any character except newline: '.'

optional instance: $R? = \epsilon \mid R$

one or more instances: $R+ == RR^*$

Equivalence of Regular Expressions and Finite Automata

Regular expressions and finite automata have equivalent expressive power:

For every regular expression R , there is a corresponding FA that accepts the set of strings generated by R .
For every FA A there is a corresponding regular expression that generates the set of strings accepted by A .

The proof is in two parts:

- i. an algorithm that, given a regular expression R , produces an FA A such that $L(A) == L(R)$.
- ii. an algorithm that, given an FA A , produces a regular expression R such that $L(R) == L(A)$.

Our construction of FA from regular expressions will allow "epsilon transitions" (a transition from one state to another with epsilon as the label). Such a transition is always possible, since epsilon (or the empty string) can be said to exist between any two input symbols. We can show that such epsilon transitions are a notational convenience; for every FA with epsilon transitions there is a corresponding FA without them.

Constructing an FA from an RE

We begin by showing how to construct an FA for the operands in a regular expression.

If the operand is a character c , then our FA has two states, s_0 (the start state) and s_F (the final, accepting state), and a transition from s_0 to s_F with label c .

If the operand is epsilon, then our FA has two states, s_0 (the start state) and s_F (the final, accepting state), and an epsilon transition from s_0 to s_F .

If the operand is null, then our FA has two states, s_0 (the start state) and s_F (the final, accepting state), and no transitions.

Given FA for R_1 and R_2 , we now show how to build an FA for R_1R_2 , $R_1|R_2$, and R_1^* . Let A (with start state a_0 and final state a_F) be the machine accepting $L(R_1)$ and B (with start state b_0 and final state b_F) be the machine accepting $L(R_2)$.

The machine C accepting $L(R_1R_2)$ includes A and B , with start state a_0 , final state b_F , and an epsilon transition from a_F to b_0 .

The machine C accepting $L(R_1|R_2)$ includes A and B , with a new start state c_0 , a new final state c_F , and epsilon transitions from c_0 to a_0 and b_0 , and from a_F and b_F to c_F .

The machine C accepting $L(R_1^*)$ includes A , with a new start state c_0 , a new final state c_F , and epsilon transitions from c_0 to a_0 and c_F , and from a_F to a_0 , and from a_F to c_F .

Eliminating Epsilon Transitions

If we can eliminate epsilon transitions from an FA, then our construction of an FA from a regular expression (which yields an FA with epsilon transitions) can be completed.

Observe that epsilon transitions are similar to nondeterminism in that they offer a choice: an epsilon transition allows us to stay in a state or move to a new state, regardless of the input symbol.

If starting in state s_1 , we can reach state s_2 via a series of epsilon transitions followed by a transition on input symbol x , we can replace all of the epsilon transitions with a single transition from s_1 to s_2 on symbol x .

Algorithm for Eliminating Epsilon Transitions

We can build a finite automaton F_2 with no epsilon transitions from a finite automaton F_1 containing epsilon transitions as follows:

The states of F_2 are all the states of F_1 that have an entering transition labeled by some symbol other than epsilon, plus the start state of F_1 , which is also the start state of F_2 .

For each state in F_1 , determine which other states are reachable via epsilon transitions only. If a state of F_1 can reach a final state in F_1 via epsilon transitions, then the corresponding state is a final state in F_2 .

For each pair of states i and j in F_2 , there is a transition from state i to state j on input x if there exists a state k that is reachable from state i via epsilon transitions in F_1 , and there is a transition in F_1 from state k to state j on input x .

Constructing an RE from an FA

To construct a regular expression from a DFA (and thereby complete the proof that regular expressions and finite automata have the same expressive power), we replace each state in the DFA one by one with a corresponding regular expression.

Just as we built a small FA for each operator and operand in a regular expression, we will now build a small regular expression for each state in the DFA.

The basic idea is to eliminate the states of the FA one by one, replacing each state with a regular expression that generates the portion of the input string that labels the transitions into and out of the state being eliminated.

Algorithm for Constructing an RE from an FA

Given a DFA F we construct a regular expression R such that $L(F) == L(R)$.

We preprocess the FA, turning the labels on transitions into regular expressions. If there is a transition with label $\{a,b\}$, then we replace the label with the regular expression $a | b$. If there is no transition from a state to itself, we can add one with the label NULL.

For each accepting state s_F in F , eliminate all states in F except the start state s_0 and s_F .

Lab 03: Design a code for regular expression to DFA and vice versa

To eliminate a state s_E , consider all pairs of states s_A and s_B such that there is a transition from s_A to s_E with label R_1 , a transition from s_E to s_E with label R_2 (possibly null, meaning no transition), and a transition from s_E to s_B with label R_3 . Introduce a transition from s_A to s_B with label $R_1R_2^*R_3$. If there is already a transition from s_A to s_B with label R_4 , then replace that label with $R_4R_1R_2^*R_3$.

After eliminating all states except s_0 and s_F :

If $s_0 == s_F$, then the resulting regular expression is R_1^* , where R is the label on the transition from s_0 to s_0 .

If $s_0 != s_F$, then assume the transition from s_0 to s_0 is labeled R_1 , the transition from s_0 to s_F is labeled R_2 , the transition from s_F to s_F is labeled R_3 , and the transition from s_F to s_0 is labeled R_4 . The resulting regular expression is $R_1^*R_2(R_3 | R_4R_1^*R_2)^*$

Let RF_i be the regular expression produced by eliminating all the states except s_0 and s_{Fi} . If there are n final states in the DFA, then the regular expression that generates the strings accepted by the original DFA is $RF_1 | RF_2 | \dots | RF_n$.

Regular Expressions Library (since C++11)

The regular expressions library provides a class that represents regular expressions, which are a kind of mini-language used to perform pattern matching within strings. Almost all operations with regexes can be characterized by operating on several of the following objects:

Target Sequence: The character sequence that is searched for a pattern. This may be a range specified by two iterators, a null-terminated character string or a `std::string`.

Pattern: This is the regular expression itself. It determines what constitutes a match. It is an object of type `std::basic_regex`, constructed from a string with special syntax. See `regex_constants::syntax_option_type` for the description of supported syntax variations.

Matched Array: The information about matches may be retrieved as an object of type `std::match_results`.

Replacement String: This is a string that determines how to replace the matches, see `regex_constants::match_flag_type` for the description of supported syntax variations.

Main Classes

These classes encapsulate a regular expression and the results of matching a regular expression within a target sequence of characters.

| | |
|------------------------------------|---|
| <code>basic_regex (C++11)</code> | regular expression object (class template) |
| <code>sub_match (C++11)</code> | identifies the sequence of characters matched by a sub-expression (class template) |
| <code>match_results (C++11)</code> | identifies one regular expression match, including all sub-expression matches (class template) |

Algorithms

These functions are used to apply the regular expression encapsulated in a regex to a target sequence of characters.

| | |
|-----------------------|---|
| regex_match (C++11) | attempts to match a regular expression to an entire character sequence (function template) |
| regex_search (C++11) | attempts to match a regular expression to any part of a character sequence (function template) |
| regex_replace (C++11) | replaces occurrences of a regular expression with formatted replacement text (function template) |

Iterators

The regex iterators are used to traverse the entire set of regular expression matches found within a sequence.

| | |
|------------------------------|---|
| regex_iterator (C++11) | iterates through all regex matches within a character sequence (class template) |
| regex_token_iterator (C++11) | iterates through the specified sub-expressions within all regex matches in a given string or through unmatched substrings (class template) |

Exceptions

This class defines the type of objects thrown as exceptions to report errors from the regular expressions library.

| | |
|---------------------|--|
| regex_error (C++11) | reports errors generated by the regular expressions library (class) |
|---------------------|--|

Traits

The regex traits class is used to encapsulate the localizable aspects of a regex.

| | |
|----------------------|---|
| regex_traits (C++11) | provides meta-information about a character type, required by the regex library (class template) |
|----------------------|---|

Constants

Defined in namespace std::regex_constants

| | |
|----------------------------|---|
| syntax_option_type (C++11) | general options controlling regex behavior (typedef) |
| match_flag_type (C++11) | options specific to matching (typedef) |
| error_type (C++11) | describes different types of matching errors (typedef) |

Lab 03: Design a code for regular expression to DFA and vice versa

Program to use `regex_search()` for finding the pattern in a given string with other parameters.

```
#include <iostream>
#include <iterator>
#include <string>
#include <regex>

int main()
{
    std::string s = "Some people, when confronted with a problem, think "
        "\"I know, I'll use regular expressions.\" "
        "Now they have two problems.";

    std::regex self_regex("REGULAR EXPRESSIONS",
        std::regex_constants::ECMAScript | std::regex_constants::icase);
    if (std::regex_search(s, self_regex)) {
        std::cout << "Text contains the phrase 'regular expressions'\n";
    }

    std::regex word_regex("(\\w+)");
    auto words_begin =
        std::sregex_iterator(s.begin(), s.end(), word_regex);
    auto words_end = std::sregex_iterator();

    std::cout << "Found "
        << std::distance(words_begin, words_end)
        << " words\n";

    const int N = 6;
    std::cout << "Words longer than " << N << " characters:\n";
    for (std::sregex_iterator i = words_begin; i != words_end; ++i) {
        std::smatch match = *i;
        std::string match_str = match.str();
        if (match_str.size() > N) {
            std::cout << " " << match_str << '\n';
        }
    }

    std::regex long_word_regex("(\\w{7,})");
    std::string new_s = std::regex_replace(s, long_word_regex, "[$&]");
    std::cout << new_s << '\n';

    return 0;
}
```

- i. Write your observations after executing this piece of code.
- ii. What is the purpose of using `[$&]`, `\\w{7, }` and `\\w+`.
- iii. What you understand by ECMAScript?

Lab 03: Design a code for regular expression to DFA and vice versa

Program to use `regex_search()` for finding the pattern in a given string.

```
// C++ program to demonstrate working of regex_search()
#include <iostream>
#include <regex>
#include <string.h>
using namespace std;

int main()
{
    // Target sequence
    string s = "Pakistan is my Love, I love Pakistan";

    // An object of regex for pattern to be searched
    regex r("Pakistan[a-zA-Z]+");

    // flag type for determining the matching behavior
    // here it is for matches on 'string' objects
    smatch m;

    // regex_search() for searching the regex pattern
    // 'r' in the string 's'. 'm' is flag for determining
    // matching behavior.
    regex_search(s, m, r);

    // for each loop
    for (auto x : m)
        cout << x << " ";

    return 0;
}
```

After executing the above code, answer the following:

- i. What is the purpose of the above regex `r`?
- ii. What if we need that this regular expression will cater space, how does the regular expression changes?
- iii. What is the purpose of `auto` in this code?
- iv. What does the function `regex_search()` will do?
- v. Sentence without space
- vi. Change the sentence according to your favorite quote and then check.

Lab 03: Design a code for regular expression to DFA and vice versa

Program to demonstrate the working of `regex_match()` for a given string.

```
// C++ program to demonstrate working of regex_match()
#include <iostream>
#include <regex>

using namespace std;
int main()
{
    string a = "Pakistan is my Love, I love Pakistan";

    // Here b is an object of regex (regular expression)
    regex b("(Pakistan)(.*)"); // Pakistan followed by any character

    // regex_match function matches string a against regex b
    if ( regex_match(a, b) )
        cout << "String 'a' matches regular expression 'b' \n";

    // regex_match function for matching a range in string
    // against regex b
    if ( regex_match(a.begin(), a.end(), b) )
        cout << "String 'a' matches with regular expression "
              << "'b' in the range from 0 to string end\n";

    return 0;
}
```

After executing the above code, answer the followings:

- i. What is the purpose of using `regex_match`?
- ii. What is the purpose of using `a.begin()` and `a.end()` in this code?
- iii. Can we change the above code? If yes enter the regex according to your criteria and show the result.

Program of using regex_replace()

```
// C++ program to demonstrate working of regex_replace()
#include <iostream>
#include <string>
#include <regex>
#include <iterator>
using namespace std;

int main()
{
    string s = "I am looking for Pakistanifamily \n";

    // matches words beginning by "Pakis"
    regex r("Pakis[a-zA-z]+");

    // regex_replace() for replacing the match with 'friend'
    cout << std::regex_replace(s, r, "friend");

    string result;

    // regex_replace( ) for replacing the match with 'friend'
    regex_replace(back_inserter(result), s.begin(), s.end(),
        r, "friend");

    cout << result;

    return 0;
}
```

Change the conditions in the given sentence and then check the outcomes. Observe the difference between them.

- i. Change the 2nd with some other word and observe the difference.
- ii. Change the sentence in terms of grammar and see if this can work.
- iii. Write the purpose of string in regex r? What other conditions we can use? Give five examples.

Lab 03: Design a code for regular expression to DFA and vice versa

Program to demonstrate the validation of incoming data.

```
//Check the below program to demonstrate how you can use the regex to validate incoming data.

#include <iostream>
#include <regex>
#include <string>
using namespace std;

int main()
{
    string input;
    regex integer_expr("(\\+|-)?[[:digit:]]+");
    //As long as the input is correct ask for another number
    while(true)
    {
        cout<<"Enter the input: ";
        cin>>input;
        if(!cin) break;
        //Exit when the user inputs q
        if(input=="q")
            break;
        if(regex_match(input,integer_expr))
            cout<<"Input is an integer"<<endl;
        else
            {cout<<"Invalid input : Not an integer"<<endl;}
    }

    return 0;
}
```

Lab 03: Design a code for regular expression to DFA and vice versa

Another C++ program to demonstrate the `regex_replace()` function.

```
#include <iostream>
#include <string>
#include <regex>
#include <iterator>
using namespace std;

int main()
{
    string mystr = "This is software testing Help portal \n";

    cout<<"Input string: "<<mystr<<endl;

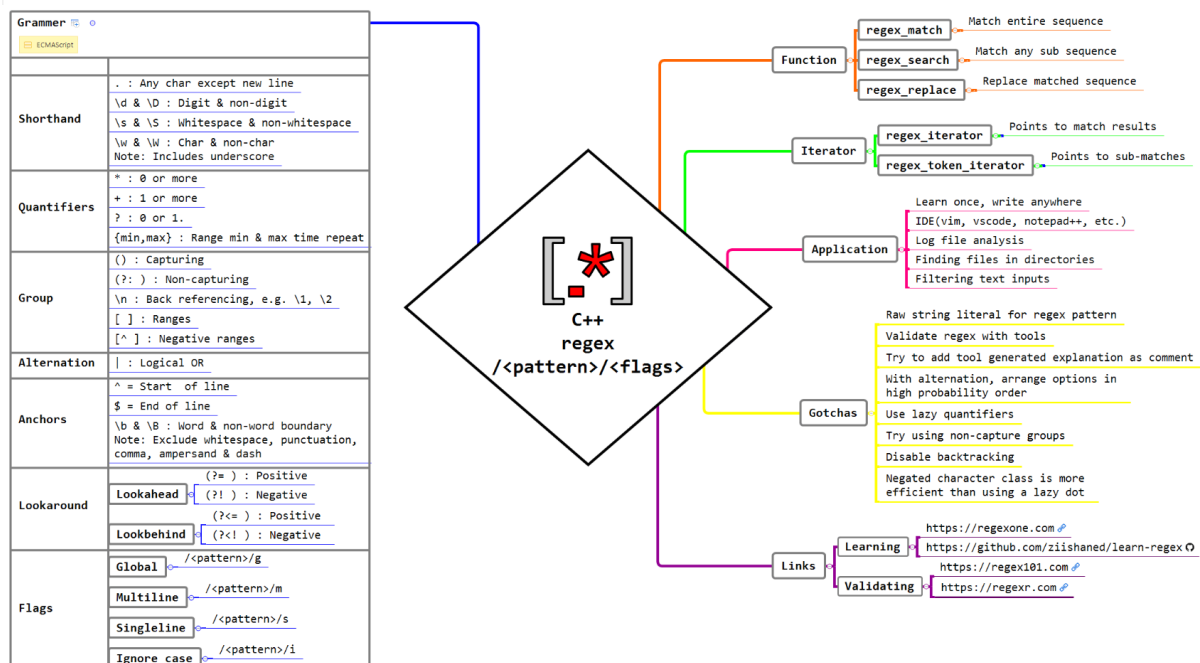
    // regex to match string beginning with 'p'
    regex regexp("p[a-zA-z]+");
    cout<<"Replace the word 'portal' with word 'website' : ";
    // regex_replace() for replacing the match with the word 'website'
    cout << regex_replace(mystr, regexp, "website");

    string result;

    cout<<"Replace the word 'website' back to 'portal': ";
    // regex_replace( ) for replacing the match back with 'portal'
    regex_replace(back_inserter(result), mystr.begin(), mystr.end(),
        regexp, "portal");

    cout << result;

    return 0;
}
```



Lab 03: Design a code for regular expression to DFA and vice versa

3. Lab tasks

Task 1

Construct the automata machines for any five of the following given conditions of regular expressions:

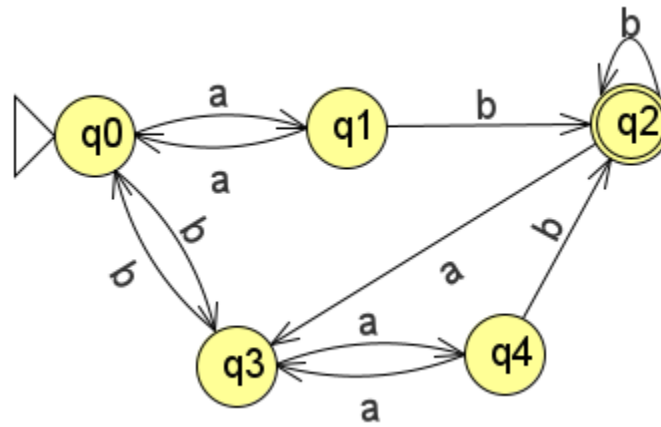
| <i>regular language</i> | <i>regular expression</i> | <i>in the language</i> | <i>not in the language</i> |
|---|---------------------------------------|---|--------------------------------------|
| binary alphabet | | | |
| <i>fifth-to-last symbol is a</i> | $(a b)^*a(a b)(a b)(a b)(a b)$ | aaaaa bbbabbbb bbbbbababababa | a bbbbbbba aaaaaaaaaaaaa |
| <i>contains the substring abba</i> | $(a b)^*abba(a b)^*$ | abba aababbabbababba bbbbbbbabababbbb | abb bbabaab aaaaaaaaaaaaa |
| <i>does not contain the substring bbb</i> | $(bba ba a^*)^*(a^* b bb)$ | aa ababababbaba aaaaaaaaaaaaab | bbb ababbbbabab bbbbbbbbb |
| <i>number of b symbols is a multiple of 3</i> | $a^* (a^*ba^*ba^*ba^*)^*$ | bbb aaa bbbaababbba | b baaaaaaab baabbbbaaaab |
| decimal digits | | | |
| <i>positive integer divisible by 5</i> | $5 (1 2 \dots 9)(0 1 \dots 9)^*(0 5)$ | 5 200 9836786785 | 1 0005 3452345234 |
| <i>positive ternary number</i> | $(1 2)(0 1 2)^*$ | 11 2210221 | 011 19 9836786785 |
| lowercase letters | | | |
| <i>contains the trigraph spb</i> | $(a b c \dots z)^*spb(a b c \dots z)$ | raspberry crispbread | subspace subspecies |
| <i>uses only the top row of the keyboard</i> | $(q w e r t y u i o p)^*$ | typewriter reporter | alfalfa paratrooper |
| genetic code | | | |
| <i>fragile X syndrome pattern</i> | $GCG(CGG AGG)^*CTG$ | GCGCTG GCGCGGCTG GCGCGGAGGCTG | GCGCGG CGGCGGCGGCTG GCGCAGGCTG |

Task 2

Write a program in C++ that can help user to enter the options using switch case to perform the operations of regular expressions based on their requirements.

Such as: match, replace, etc.

a.



b.

