

# CS143 Spring 2022 – Written Assignment 1

Thursday, April 14, 2022 11:59 PM PDT

This assignment covers regular languages, finite automata, and lexical analysis. You may discuss this assignment with other students and work on the problems together. However, your write-up should be your own individual work, and you should indicate in your submission who you worked with, if applicable. Assignments can be submitted electronically through Gradescope as a PDF by 11:59 PM PDT. Please review the course policies for more information: <https://web.stanford.edu/class/cs143/policies/>. A L<sup>A</sup>T<sub>E</sub>X template for writing your solutions is available on the course website. To create finite automata diagrams, you can either use the TikZ package directly by following the examples in the template, or a tool like <https://madebyevan.com/fsm/>.

1. Write regular expressions for the following languages over the alphabet  $\Sigma = \{0, 1\}$ . Hint: some of these languages may include  $\varepsilon$ .

- (a) The set of all strings whose 1<sup>st</sup>, 3<sup>rd</sup>, 5<sup>th</sup>,  $\dots$ , characters are the same.
- (b) The set of all strings that represent the concatenation of one odd number, one even number, and another odd number expressed in binary. (E.g., 01 10 01, but not 0110.)
- (c) The set of all strings, except the string 0000.

2. Draw DFAs for each of the languages from question 1. Note that a DFA must have a transition defined for every state and symbol pair. You must take this fact into account for your transformations. Your DFAs should not have more than 10 states.

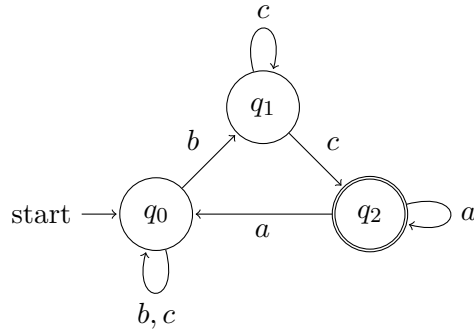
Notice that a short regular expression does not automatically imply a DFA with few states, nor vice versa.

3. Using the techniques covered in class, transform the following NFAs over the alphabet  $\{a, b, c\}$  into DFAs. Your DFAs should not have more than 10 states. Note that a DFA must have a transition defined for every state and symbol pair, whereas a NFA need not. You must take this fact into account for your transformations. Hint: Is there a subset of states the NFA transitions to when fed a symbol for which the set of current states has no explicit transition?

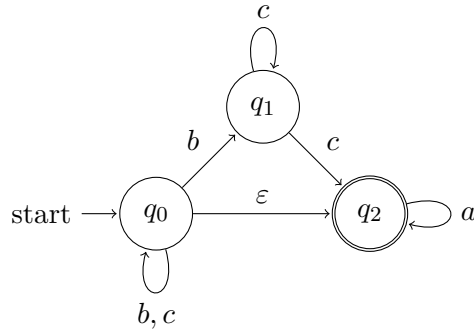
Also include a mapping from each state of your DFA to the corresponding states of the original NFA. Specifically, a state  $s$  of your DFA maps to the set of states  $Q$  of the NFA such that an input string stops at  $s$  in the DFA if and only if it stops at one of the states in  $Q$  in the NFA.

Tip: for readability, states in the DFA may be labeled according to the set of states they represent in the NFA. For example, state  $q_{012}$  in the DFA would correspond to the set of states  $\{q_0, q_1, q_2\}$  in the NFA, whereas state  $q_{13}$  would correspond to set of states  $\{q_1, q_3\}$  in the NFA.

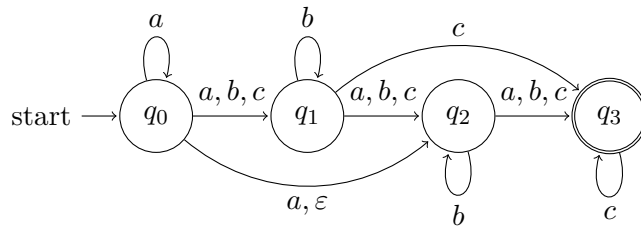
(a)



(b)



(c)



4. Let  $L$  be a language over  $\Sigma = \{a, b, c\}$ , such that string  $w$  is in  $L$  if and only if  $w$  is not  $\varepsilon$  and the last character of  $w$  appears at most twice in  $w$ .

Examples of strings in  $L$ : aa, aba, babababc.

Examples of strings **not** in  $L$ :  $\varepsilon$ , bbb, cabaa

Draw an NFA for  $L$ . Your solution should have no more than 15 states.

5. Consider the following tokens and their associated regular expressions, given as a **flex** scanner specification:

```
%%
01?1                                printf("apple");
0(10)+10                            printf("banana");
(1011*0|0100*1)                    printf("coconut");
```

Give an input to this scanner such that the output string is  $(\text{apple})^3((\text{banana})^2 \text{ coconut})^2$ , where  $A^i$  denotes  $A$  repeated  $i$  times. (And, of course, the parentheses are not part of the output.) You may use similar shorthand notation in your answer.

6. Recall from the lecture that, when using regular expressions to scan an input, we resolve conflicts by taking the largest possible match at any point. That is, if we have the following **flex** scanner specification:

```
%%  
do                { return T_Do; }  
[A-Za-z_][A-Za-z0-9_]* { return T_Identifier; }
```

and we see the input string “**dot**”, we will match the second rule and emit `T_Identifier` for the whole string, not `T_Do`.

However, it is possible to have a set of regular expressions for which we can tokenize a particular string, but for which taking the largest possible match will fail to break the input into tokens. Give an example of no more than two regular expressions and an input string such that: a) the string can be broken into substrings, where each substring matches one of the regular expressions, b) our usual lexer algorithm, taking the largest match at every step, will fail to break the string in a way in which each piece matches one of the regular expressions. Explain how the string can be tokenized and why taking the largest match won't work in this case.

As a challenge (not necessary for credit), try to find a solution that only uses one regular expression.