

CS/ECE 374 A (Spring 2020)

Old HW2 Problems with Solutions

Problem OLD.2.1: *C comments* are the set of strings over alphabet $\Sigma = \{*, /, A, \square, \langle\langle\text{Enter}\rangle\rangle\}$ that form a proper comment in the C program language and its descendants, like C++ and Java. Here $\langle\langle\text{Enter}\rangle\rangle$ represents the newline character, \square represents any other whitespace character (like the space and tab characters), and A represents any non-whitespace character other than $*$ or $/$.¹ There are two types of C comments:

- Line comments: Strings of the form $// \dots \langle\langle\text{Enter}\rangle\rangle$.
- Block comments: Strings of the form $/* \dots */$.

Following the C99 standard, we explicitly disallow *nesting* comments of the same type. A line comment starts with $//$ and ends at the first $\langle\langle\text{Enter}\rangle\rangle$ after the opening $//$. A block comment starts with $/*$ and ends at the first $*/$ completely after the opening $/*$; in particular, every block comment has at least two $*$ s. For example, each of the following strings is a valid C comment:

- $/* */$
- $//\square//\square \langle\langle\text{Enter}\rangle\rangle$
- $/* //\square * \square \langle\langle\text{Enter}\rangle\rangle * */$
- $/* \square //\square \langle\langle\text{Enter}\rangle\rangle \square * /$

On the other hand, *none* of the following strings is a valid C comments:

- $/* /$
- $//\square//\square \langle\langle\text{Enter}\rangle\rangle \square \langle\langle\text{Enter}\rangle\rangle$
- $/* \square / * \square * / \square * /$

- (a) Describe a DFA that accepts the set of all C comments.
- (b) Describe a DFA that accepts the set of all strings composed entirely of blanks (\square), newlines ($\langle\langle\text{Enter}\rangle\rangle$), and C comments.

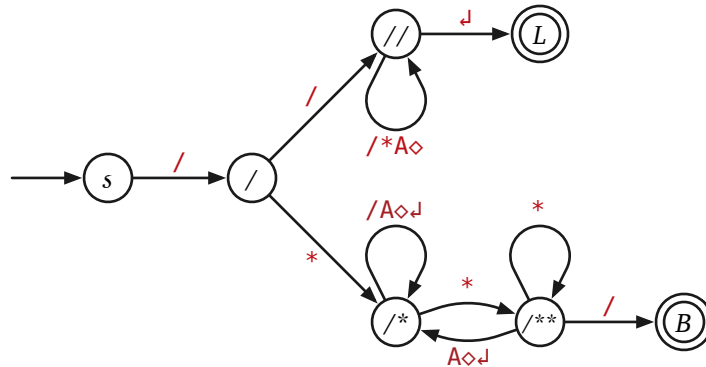
You must explain in English how your DFAs work. Drawings or formal descriptions without English explanations will receive no credit, even if they are correct.

¹The actual C commenting syntax is considerably more complex than described here, because of character and string literals.

- The opening $/*$ or $//$ of a comment must not be inside a string literal ($" \dots "$) or a (multi-)character literal ($' \dots '$).
- The opening double-quote of a string literal must not be inside a character literal ($" "$) or a comment.
- The closing double-quote of a string literal must not be escaped ($\backslash "$)
- The opening single-quote of a character literal must not be inside a string literal ($" \dots ' \dots "$) or a comment.
- The closing single-quote of a character literal must not be escaped ($\backslash '$)
- A backslash escapes the next symbol if and only if it is not itself escaped ($\backslash \backslash$) or inside a comment.

Solution:

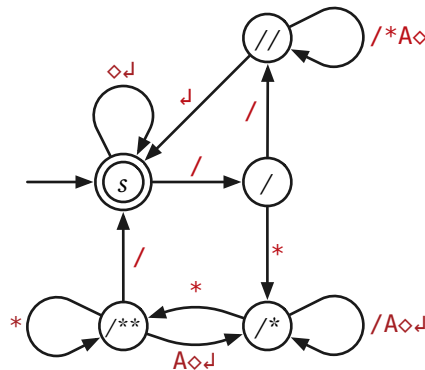
- (a) The following eight-state DFA recognizes the language of C comments. All missing transitions lead to a hidden reject state.



The states are labeled mnemonically as follows:

- *s* - We have not read anything.
- */* - We just read the initial */*.
- *//* - We are reading a line comment.
- *L* - We have read a complete line comment.
- */** - We are reading a block comment, and we did not just read a *** after the opening */**.
- */*** - We are reading a block comment, and we just read a *** after the opening */**.
- *B* - We have read a complete block comment.

- (b) By merging the accepting states of the previous DFA with the start state and adding white-space transitions at the start state, we obtain the following six-state DFA. Again, all missing transitions lead to a hidden reject state.



For example, the string `"/ * \\\" * /\" * / * \" / * \" * /` is a valid string literal (representing the 5-character string `/* \" * /`, which is itself a valid block comment!) followed immediately by a valid block comment. *For this homework question, just pretend that the characters ' , " , and \ don't exist.*

Commenting in C++ is even more complicated, thanks to the addition of *raw* string literals. Don't ask.

Some C and C++ compilers do support nested block comments, in violation of the language specification. A few other languages, like OCaml, explicitly allow nesting block comments.

The states are labeled mnemonically as follows:

- s - We are between comments.
- $/$ - We just read the initial $/$ of a comment.
- $//$ - We are reading a line comment.
- $/*$ - We are reading a block comment, and we did not just read a $*$ after the opening $/*$.
- $/**$ - We are reading a block comment, and we just read a $*$ after the opening $/*$.

Note: Generally, a DFA can be described in several ways:

- **Drawings:** Use an arrow from nowhere to indicate s , and doubled circles to indicate accepting states A . If $A = \emptyset$, say so explicitly. If your drawing omits a reject state, say so explicitly. **Draw neatly!** If we can't read your solution, we can't give you credit for it.
- **Text descriptions:** You can describe the transition function either using a 2d array, using mathematical notation, or using an algorithm.
- **Product constructions:** You must give a complete description of the states and transition functions of the DFAs you are combining (as either drawings or text), together with the accepting states of the product DFA.

In homeworks, we should correctly explain the purpose of each state *in English*. This is how you justify that your DFA is correct.

- **Deadly Sin:** ("Declare your variables.") No credit for the problem if the English description is missing, *even if the DFA is correct*. (For product constructions, explaining the states in the factor DFAs is enough.)
- DFA drawings with too many states may be penalized.
- Describing an NFA when the problem asks for a DFA will be severely penalized.

Problem OLD.2.2. Give a regular expression that describes the following language over the alphabet $\{0, 1\}$, and briefly argue why your expression is correct.

- All strings in which every nonempty maximal substring of consecutive 0s is of length 1. For instance 1001 is not in the language while 10111 is.

Solution: Let's start with some strings that are in the language: 11^*0 . In fact, the star closure of this language is also in the language: $(11^*0)^*$. We can of course suffix this with a run of ones: $(11^*0)^*1^*$. This indeed captures all the strings in the language that starts with 1, since one can break the string after each appearance of 0, which give the above expression.

So, we only have to handle, in addition, the strings in the language that starts with 0. Such a strings can be generated only by $0(11^*0)^*1^*$. As such, the overall regular expression is

$$(\varepsilon + 0)(11^*0)^*1^*.$$

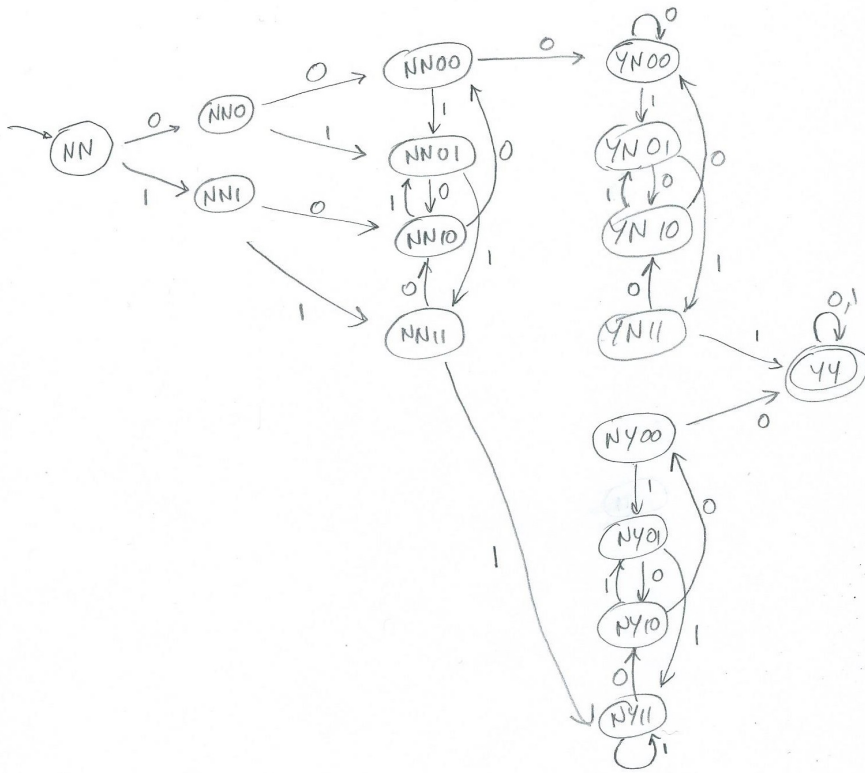
Let us prove that this is correct. Given a string w in the language, break it into strings s_1, \dots, s_k , where we cut the string w after each appearance of 0 (i.e., $w = s_1 s_2 \dots s_k$). If $s_1 = 0$, then we choose 0 in the left side, otherwise, we choose ε in the above regular expression. Similarly, if s_k contains no 0 s, then we assign it to the suffix 1^* . All other strings are assigned to the generating expression 11^*0 . To see the correctness, observe that every s_i is a non-empty run of 1 s followed by a zero.

Note: More regular expression examples can be found in lab 1b.

Problem OLD.2.3: Describe a DFA that over the alphabet $\Sigma = \{0, 1\}$ that accepts the following language:

- All strings that contain the substrings 000 and 111 .

Solution: The idea is to keep track of the last two characters we have seen, and whether we have seen 000 or 111 yet. A drawing of the DFA is given below:



- The state NNx means that we have not seen 000 nor 111 , and the last two characters we have seen are x .
- The state YNx means that we have seen 000 but have not seen 111 , and the last two characters we have seen are x .

- The state NYx means that we have not seen 000 but have seen 111 , and the last two characters we have seen are x .
- The state YY means that we have seen both 000 and 111 .

[Note: an alternative solution can be obtained by a product construction for two DFAs with 4 states each (for the language of all strings containing 000 , and the language of all strings containing 111), yielding a total of 16 states.]

Note: More DFA examples can be found in lab 2a and lab 2b.

Drawing editors: Instead of drawing your DFAs by hand, you may alternatively use a general-purpose drawing editor (e.g., I recommend `ipe`), or some drawing editor specifically for directed graphs or automata (e.g., “dot” from `graphviz`).