

**CS3342 –  
Assignment 1  
due Feb. 8, 2024  
2-day no-penalty extension until: Feb. 10,  
11:59pm**

1. (10pt) Consider a programming language that has comments as follows:

- comments start with a line beginning with three consecutive equal signs, "===", followed by any characters until the end of the line (and a newline '\n' at the end)
- comments end with a line consisting of equal signs only, at least three (and a newline '\n' at the end)

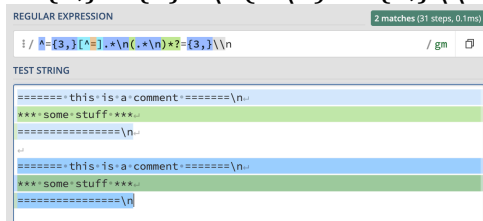
Example of a comment (newlines are shown for clarity):

```
===== this is a comment =====\n
*** some stuff ***
===== \n
```

Write a regular expression that describes precisely the above described comments. You can use the notation  $\Sigma - \{a\}$  to describe the set of all characters different from  $a$ .

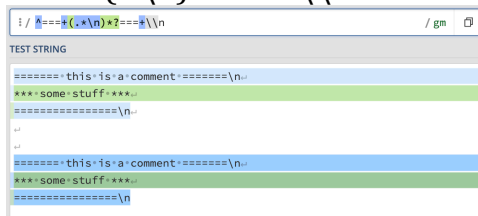
The regular expression would be:

$^{\wedge}=\{3,\}\Sigma-\{=\}.*\backslash n(.*\backslash n)^*?=\{3,\}\backslash \backslash n$



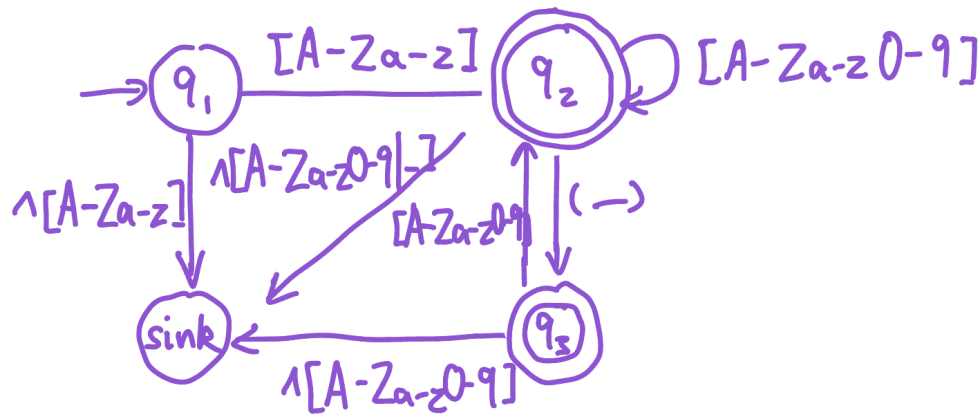
Or

$^{\wedge}===+(.*\backslash n)^*?===+\backslash \backslash n$



Reference: works on [regex101.com](https://regex101.com)

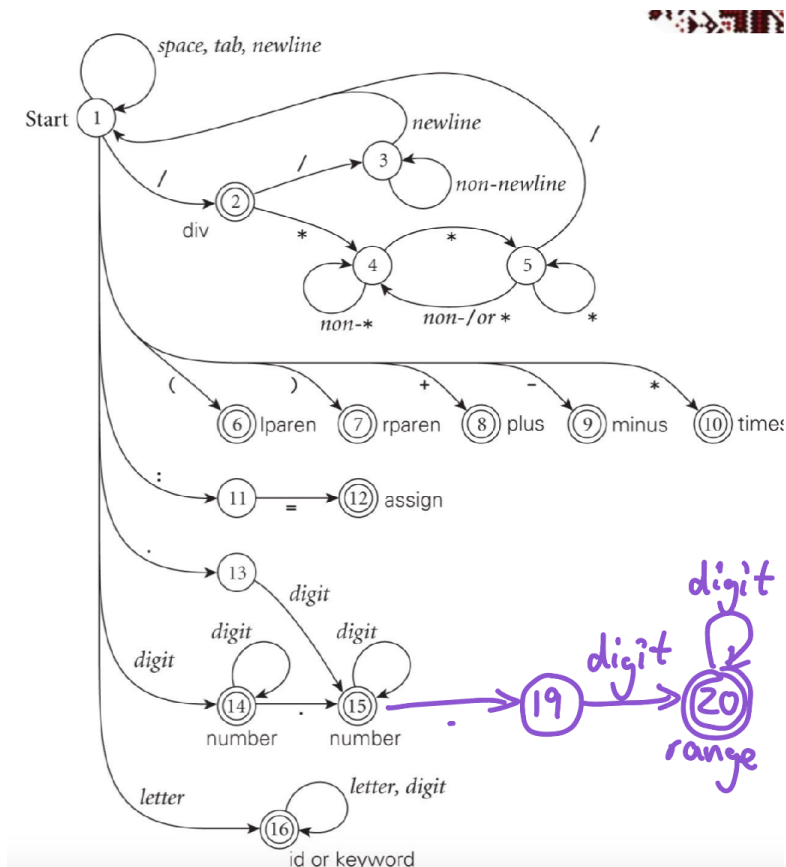
2. (5pt) A scanner is built for a language where the identifiers start with a letter followed by any number of letters, digits, or underscores, such that there are no consecutive underscores. Draw a deterministic finite automaton (DFA) with fewest states that accepts all identifiers and nothing else.



3. (10pt) Consider the simple calculator language (slide 16 of scanning) and add to it ranges:

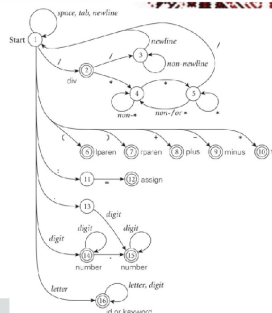
$range \rightarrow digit\ digit^* \dots digit\ digit^*$

Modify the scanner table and DFA on slide 28 to reflect this change. (Note that the DFA already requires some small changes to match the table, as discussed in class.) To save time, you are allowed to cut-and-paste the table and DFA (as figures) and only indicate the changes.



# Scanning

- Scanner table used by previous code
  - state 17: white spaces; state 18: comments
  - scan\_tab: entire table but last column
  - token\_tab: last column
  - keyword\_tab = {read, write}



State	Current input character														
	space, tab	newline	/	*	(	)	+	-	:	=	.	digit	letter	other	
1	17	17	2	10	6	7	8	9	11	-	13	14	16	-	
2	-	-	3	4	-	-	-	-	-	-	-	-	-	-	div
3	3	18	3	3	3	3	3	3	3	3	3	3	3	3	
4	4	4	4	5	4	4	4	4	4	4	4	4	4	4	
5	4	4	18	5	4	4	4	4	4	4	4	4	4	4	
6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	lparen
7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	rparen
8	-	-	-	-	-	-	-	-	-	-	-	-	-	-	plus
9	-	-	-	-	-	-	-	-	-	-	-	-	-	-	minus
10	-	-	-	-	-	-	-	-	-	-	-	-	-	-	times
11	-	-	-	-	-	-	-	-	12	-	-	-	-	-	
12	-	-	-	-	-	-	-	-	-	-	-	-	-	-	assign
13	-	-	-	-	-	-	-	-	-	-	-	15	-	-	
14	-	-	-	-	-	-	-	-	-	-	15	14	-	-	number
15	-	-	-	-	-	-	-	-	-	-	19	15	-	-	number
16	-	-	-	-	-	-	-	-	-	-	-	16	16	-	identifier
17	17	17	-	-	-	-	-	-	-	-	-	-	-	-	white.space
18	-	-	-	-	-	-	-	-	-	-	-	-	-	-	comment

19  
20

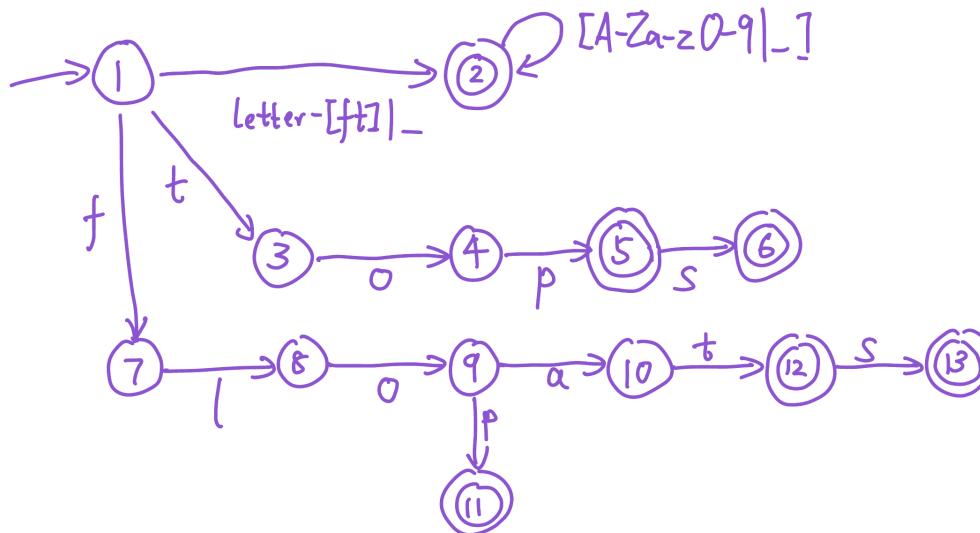
20  
20

range

4. (10pt) Keywords fit the definition of identifiers and scanners identify them as such, and then look them up a table of keywords, because otherwise the DFA is unnecessarily larger. This question addresses the size of this DFA. Assume that the identifiers and the (made-up) keywords are defined as follows:

$identifier \rightarrow ( \_ | letter )( \_ | letter | digit )^*$   
 $letter \rightarrow a | b | \dots | z$   
 $digit \rightarrow 0 | 1 | \dots | 9$   
 $keyword \rightarrow float | floats | flop | top | tops$

Draw a DFA with fewest states that recognizes the above keywords directly. That is, there will be separate accepting states for keywords and different ones for identifiers. You can use any meaningful notation such as  $letter - \{a, b, \dots\}$  or  $\neq a$  to denote various subsets of characters.



5. (15pt) Consider the following grammar,  $G$ :

- |  |  |
|--|--|
| 1. $P \rightarrow S \$ \$$                               | 6. $O \rightarrow \text{else } T$        |
| 2. $S \rightarrow \text{if (e) then } M O$               | 7. $O \rightarrow \epsilon$              |
| 3. $S \rightarrow \text{other}$                          | 8. $T \rightarrow \text{if (e) then } T$ |
| 4. $M \rightarrow \text{if (e) then } M \text{ else } M$ | 9. $T \rightarrow \text{other}$          |
| 5. $M \rightarrow \text{other}$                          |  |

- (a) (5pt) Compute the sets  $\text{FIRST}(X)$  and  $\text{FOLLOW}(X)$ , for all non-terminals  $X$ , and the sets  $\text{PREDICT}(p)$  for all production rules  $p$ .
- (b) (5pt) Prove that  $G$  an  $LL(1)$  grammar.
- (c) (5pt) We said that there is no top-down grammar for **if...then...else** statements. Does  $G$  contradict this? Prove your answer.

X	FIRST(x)	FOLLOW(x)
P	if, other	$\emptyset$
S	if, other	$\$ \$$
M	if, other	else, $\epsilon$
O	else, $\epsilon$	$\emptyset$
T	if, other	$\emptyset$

i	PREDICT(i)
1. $P \rightarrow S \$ \$$	if, other
2. $S \rightarrow \text{if (e) then } M O$	if
3. $S \rightarrow \text{other}$	other
4. $M \rightarrow \text{if (e) then } M \text{ else } M$	if
5. $M \rightarrow \text{other}$	other
6. $O \rightarrow \text{else } T$	else

- |   |       |
|---|-------|
| 7. $O \rightarrow \epsilon$                       | \$\$  |
| 8. $T \rightarrow \text{if } (e) \text{ then } T$ | if    |
| 9. $T \rightarrow \text{other}$                   | other |

Q5: LL(1) definition is on slide 17 of LL-parsing. See also slide 19 of LR-parsing.

2) Prove that G an LL(1) grammar.

For any production:  $A \rightarrow u|v$

$\text{First}(u) \cap \text{First}(v) = \emptyset$

If  $v \Rightarrow \epsilon$  then  $\text{First}(u) \cap \text{Follow}(A) = \emptyset$

At most one of u and v can derive the empty string  $\epsilon$

According to the predict table, since there is no duplicated non-terminal derive the same terminal, which mean LL(1) is hold for this grammar.

6. (50pt) Write a Python program `comm_rm.py` to remove all comments from a Python program. The program should work as follows:

Some ideas I had in mind about this question:

- the comments to be removed include single line comments and triple-quoted strings (single or double quotes) acting as comments
- if used in expressions, triple-quoted strings are not to be removed, as they do not act as comments
- you need to consider also:
  - backslash-escaped characters
  - backslash as line continuation character
  - the interaction between comments and single-quoted strings

```
> python comm_rm.py input_file.py
```

where `input_file.py` is any (correct) Python program. The program with comments removed will be saved in the file `input_file_rm.py`.

You are *not allowed* to use the regular expression capabilities of Python, as that would defeat the purpose of the question.

**READ ME! Submit all your responses as a single pdf file.** Solutions should be typed but high-quality hand-written solutions are acceptable. **Source code, if required, is submitted separately.**

**JFLAP:** You are allowed to use JFLAP to help you solve the assignment. You still need to explain clearly your solution. Also, make sure you understand what it does; JFLAP will not be available during exams!

**LLMs:** You are allowed to use LLMs (Large Language Models), such as ChatGPT, but, again, they will not be available during exams.

**L<sup>A</sup>T<sub>E</sub>X:** For those interested, the best program for scientific writing is L<sup>A</sup>T<sub>E</sub>X. It is far superior to all the other programs, it is free, and you can start using it in minutes; here is an introduction: <https://tobi.oetiker.ch/lshort/lshort.pdf>. It is also available online at <https://www.overleaf.com/>.