

Section 1

I did not change the lexical specifications which means it is the same as the problem description. I have decided to add everything under one regular expression which gives me this:

<u>id</u>	::=	Letter alphanum*
alphanum	::=	letter digit _
<u>integer</u>	::=	nonzero digit* 0
<u>float</u>	::=	integer fraction [e[+ -] integer]
fraction	::=	.digit* nonzero .0
letter	::=	a..z A..Z
digit	::=	0..9
nonzero	::=	1..9
<u>eq</u>	::=	==
<u>plus</u>	::=	+
<u>or</u>	::=	or
<u>openpar</u>	::=	(
<u>semi</u>	::=	;
<u>while</u>	::=	while
<u>localvar</u>	::=	localvar
<u>noteq</u>	::=	<>
<u>minus</u>	::=	-
<u>and</u>	::=	and
<u>closepar</u>	::=)
<u>comma</u>	::=	,
<u>if</u>	::=	if
<u>constructo r</u>	::=	constructor
<u>lt</u>	::=	>
<u>mult</u>	::=	*
<u>not</u>	::=	not
<u>opencubr</u>	::=	{
<u>dot</u>	::=	.
<u>void</u>	::=	void

<u>then</u>	::=	then
<u>attribute</u>	::=	attribute
<u>gt</u>	::=	>
<u>div</u>	::=	/
<u>closecubr</u>	::=	}
<u>colon</u>	::=	:
<u>class</u>	::=	class
<u>else</u>	::=	else
<u>function</u>	::=	function
<u>leq</u>	::=	<=
<u>assign</u>	::=	=
<u>opensqbr</u>	::=	[
<u>returntype</u>	::=	=>
<u>self</u>	::=	self
<u>read</u>	::=	read
<u>public</u>	::=	public
<u>geq</u>	::=	>=
<u>closesqbr</u>	::=]
<u>scopeop</u>	::=	::
<u>isa</u>	::=	isa
<u>write</u>	::=	write
<u>private</u>	::=	private
<u>return</u>	::=	return
<u>intnum</u>	::=	integer
<u>floatnum</u>	::=	float
<u>inlinecmt</u>	::=	// nonbreak
<u>blockcmt</u>	::=	/* anything */

**note: nonbreak represents all characters except the break line character

**note: anything represents all of the alphabet

Section 2

I have made a finite state machine for the high level overview of the lexical analyzer and also for all the subparts of it. I can't fit everything under one finite state machine.

Finite state machine for lexical analyzer:

d represents an id

i represents an integer

f represents a float

p represents a +

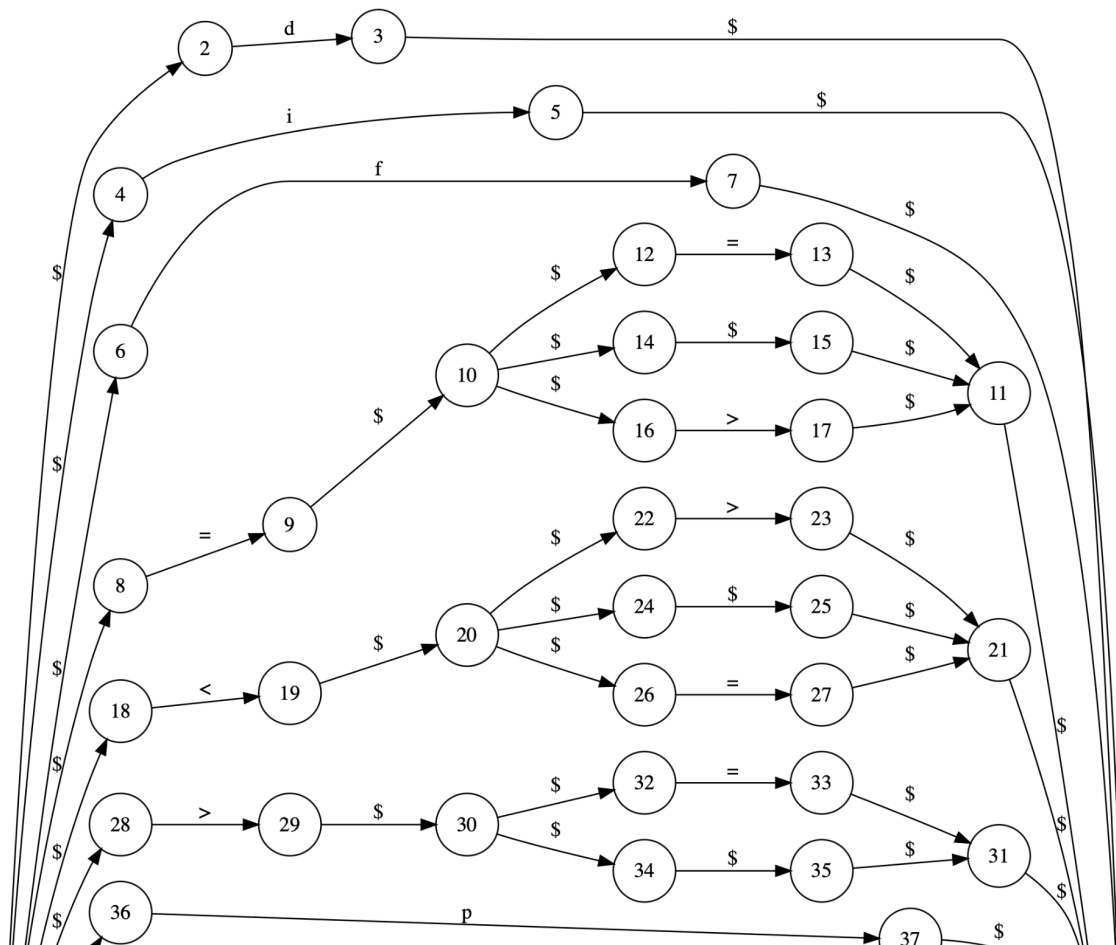
s represents a *

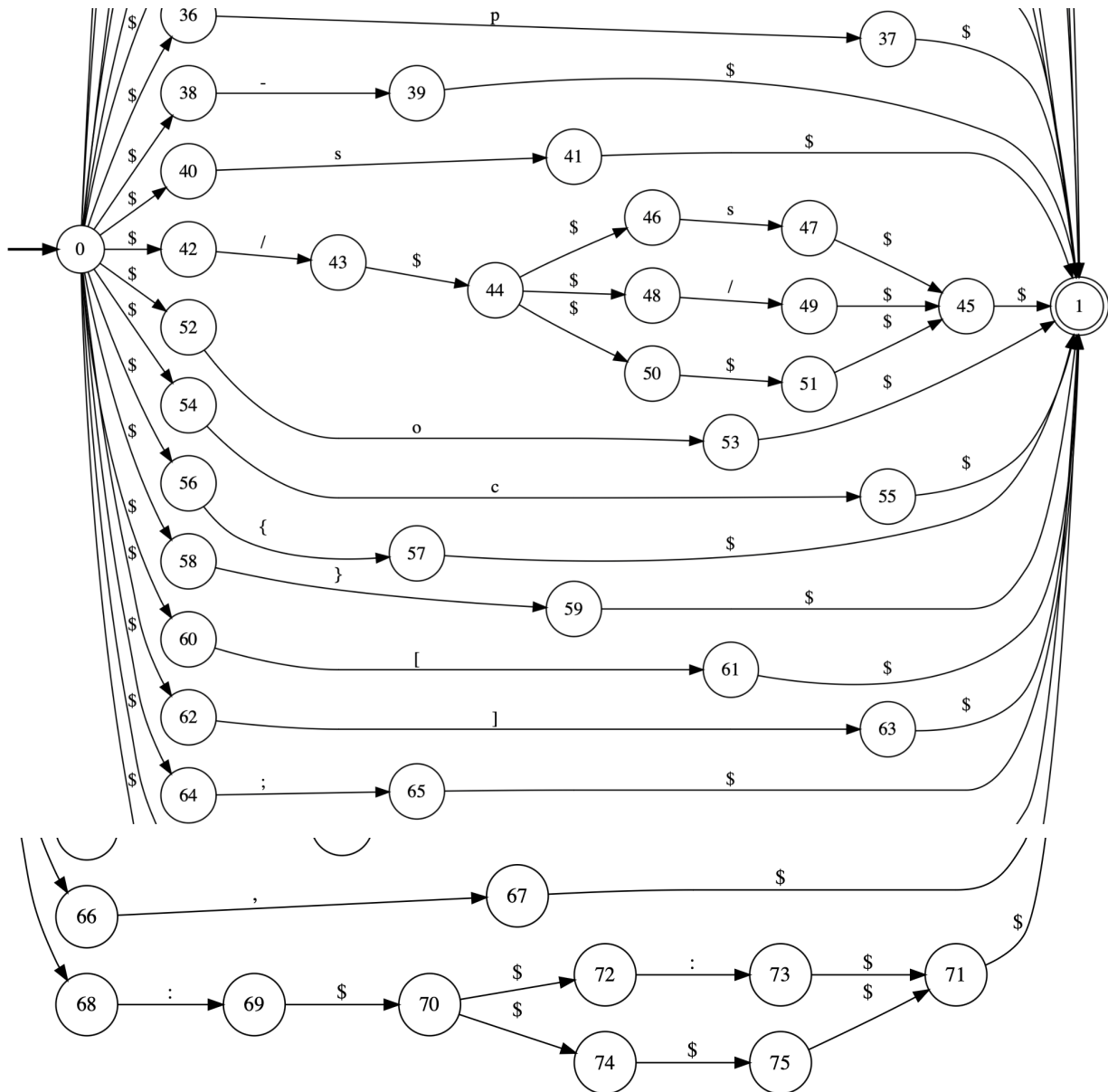
o represents a (

c represents a)

This finite state machine shows all the different paths the program can take.

In order to represent id, integer and float, I will make another finite state machine for them specifically. The diagram would be too big to show if I did it here. But we can simply take the diagrams for id, float and integer and replace them in this one. Note that I did not go in details for how comments work, but I will in later finite state machines for the comments.

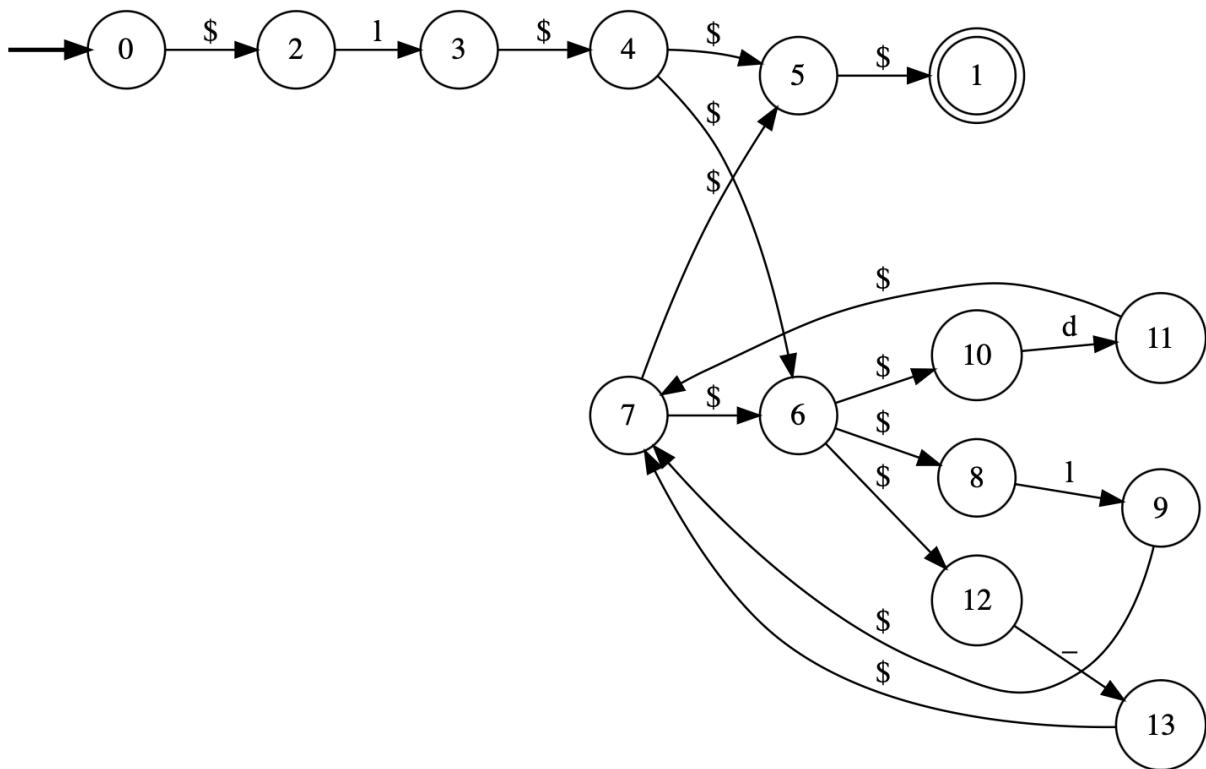




Regular expression for lexical analyzer

$d+i+f+=(+ \$ >)+<(>+ \$ +=)+>(+ \$)+p+-+s+ /+o+c+\{+\}+[+]+;+,+:(\cdot+ \$)$

Finite state machine for id
 l represents a letter
 d represents a digit

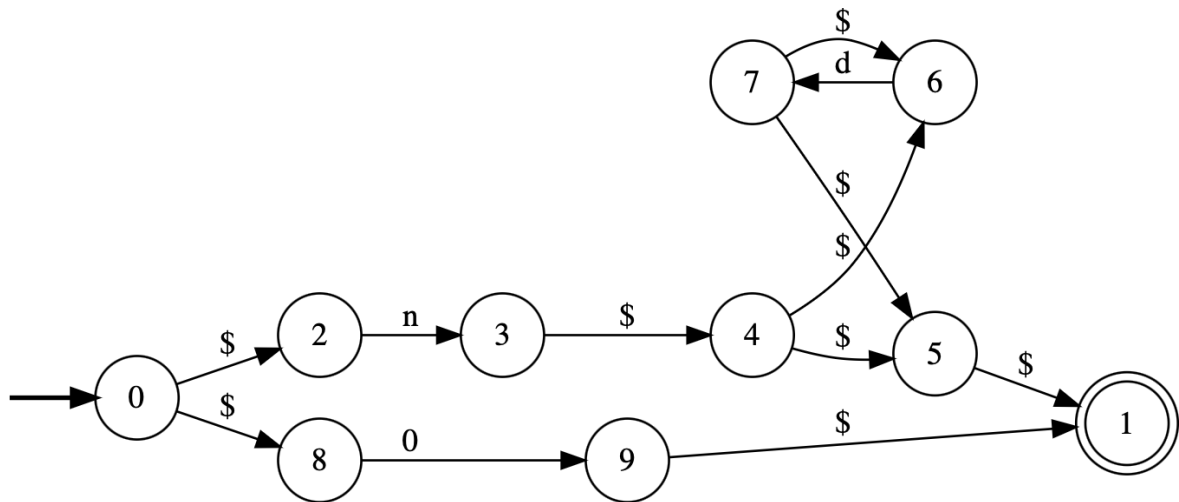


Regular expression for id: $l(l|d|_)^*$

Finite state machine for integer

N represents a nonzero

D represents a digit



Regular expression for integer: nd^*+0

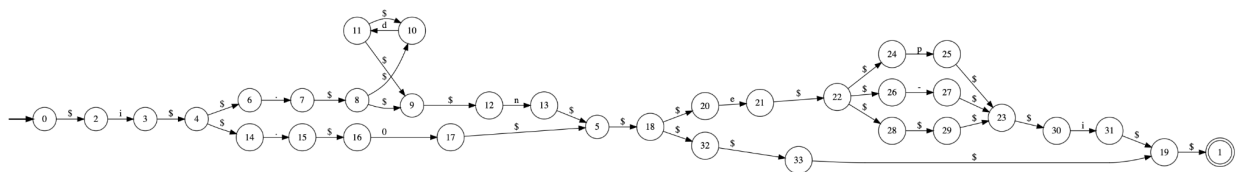
Finite state machine for float

i represents an integer

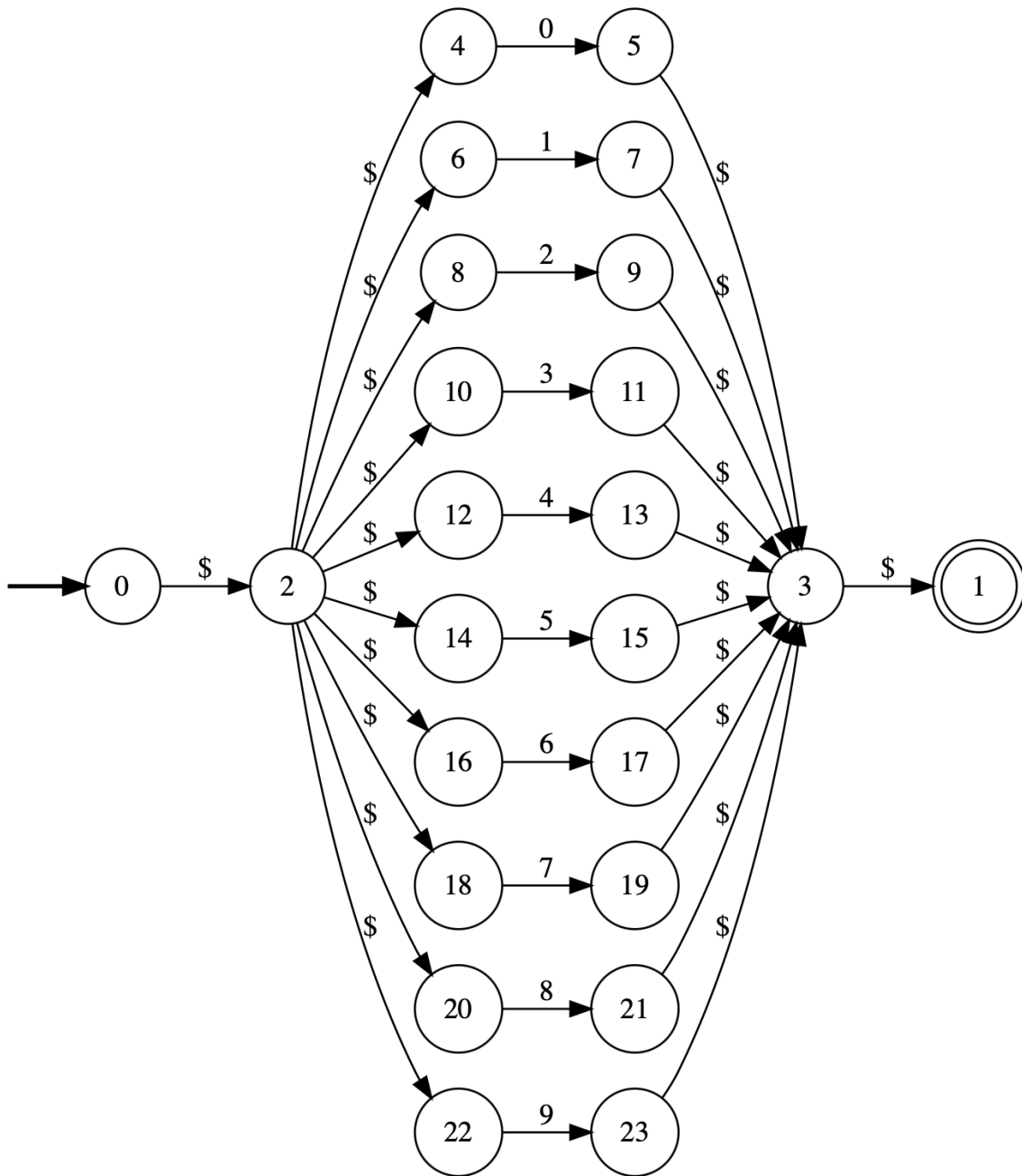
d represents a digit

n represents a nonzero

p represents a +

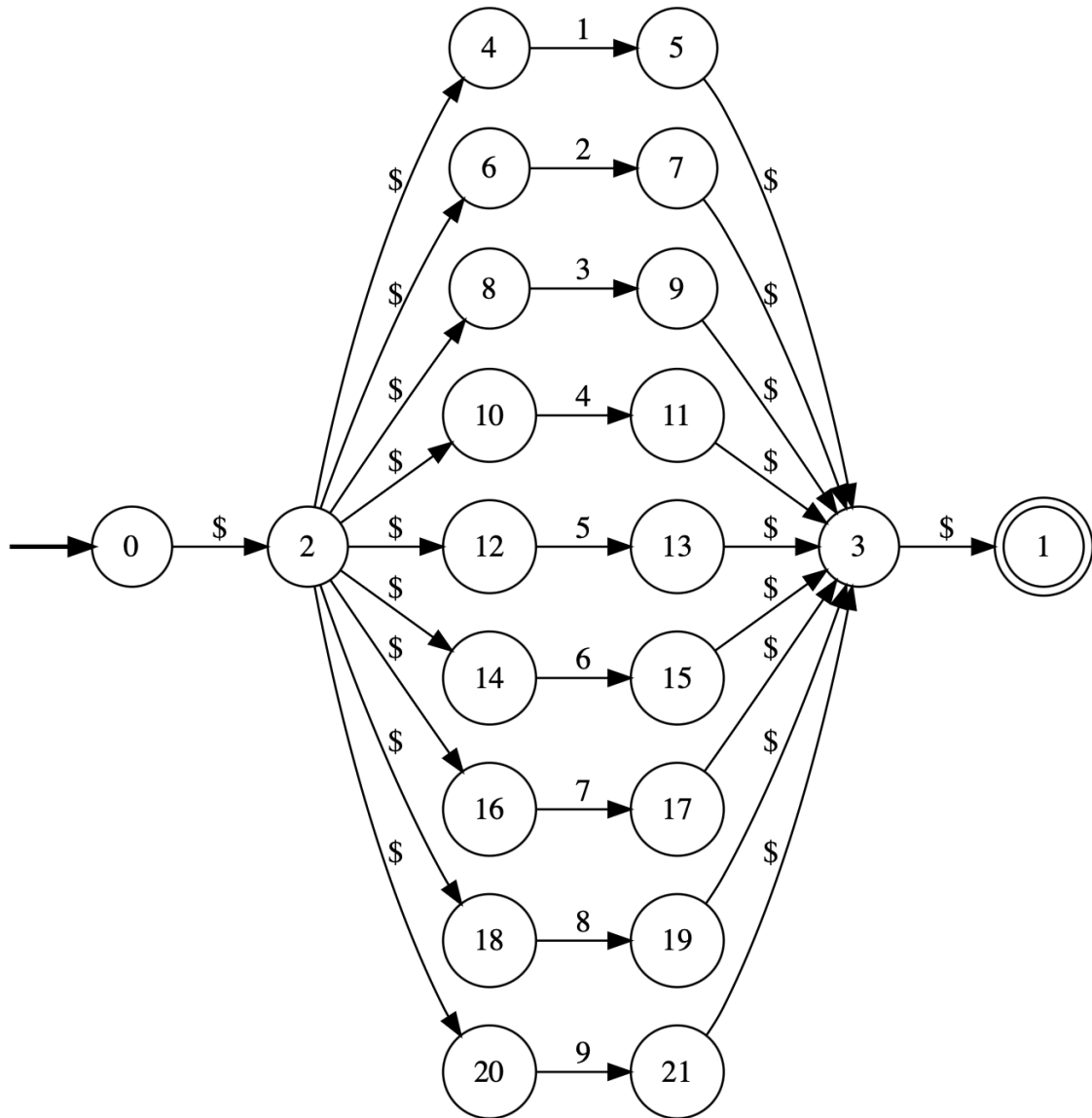


Regular expression for float: $i(d^*n+.0)(e(p+-\$)i+ \$)$
Finite state machine for digit



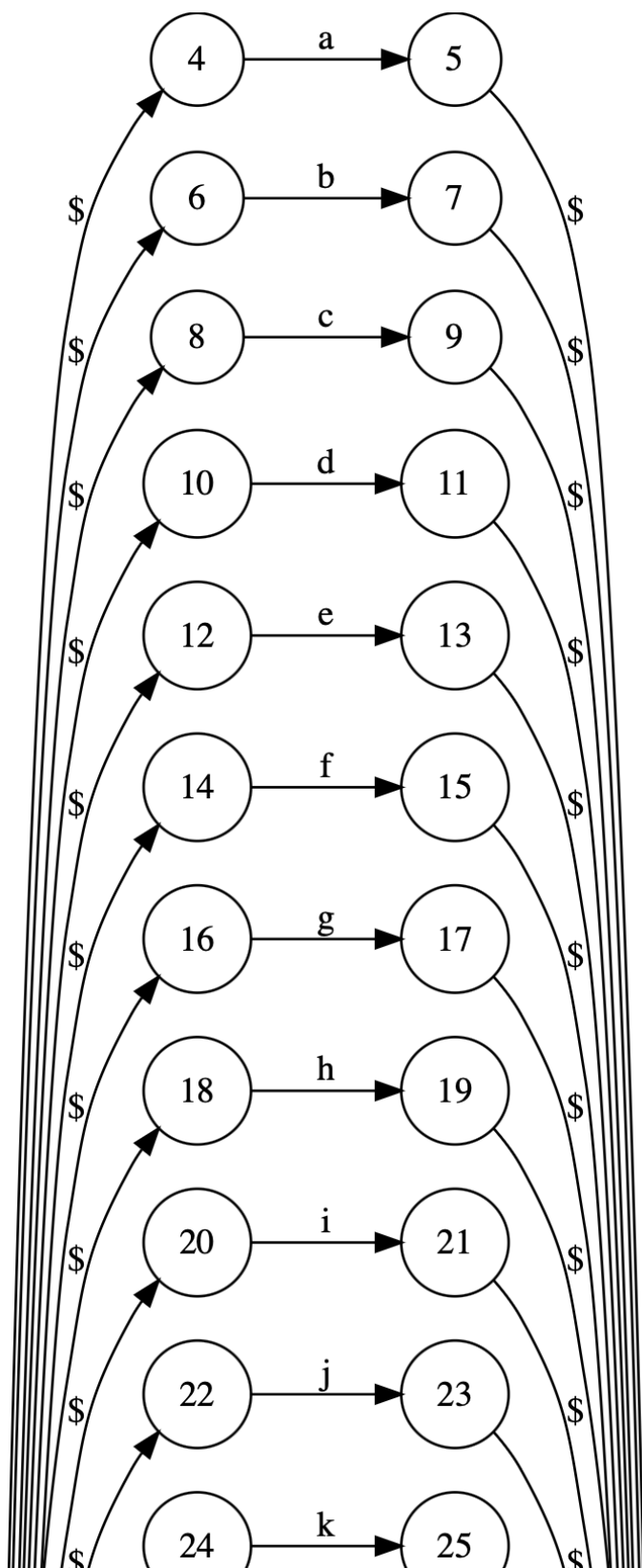
Regular expression for digit: $(1+2+3+4+5+6+7+8+9)$

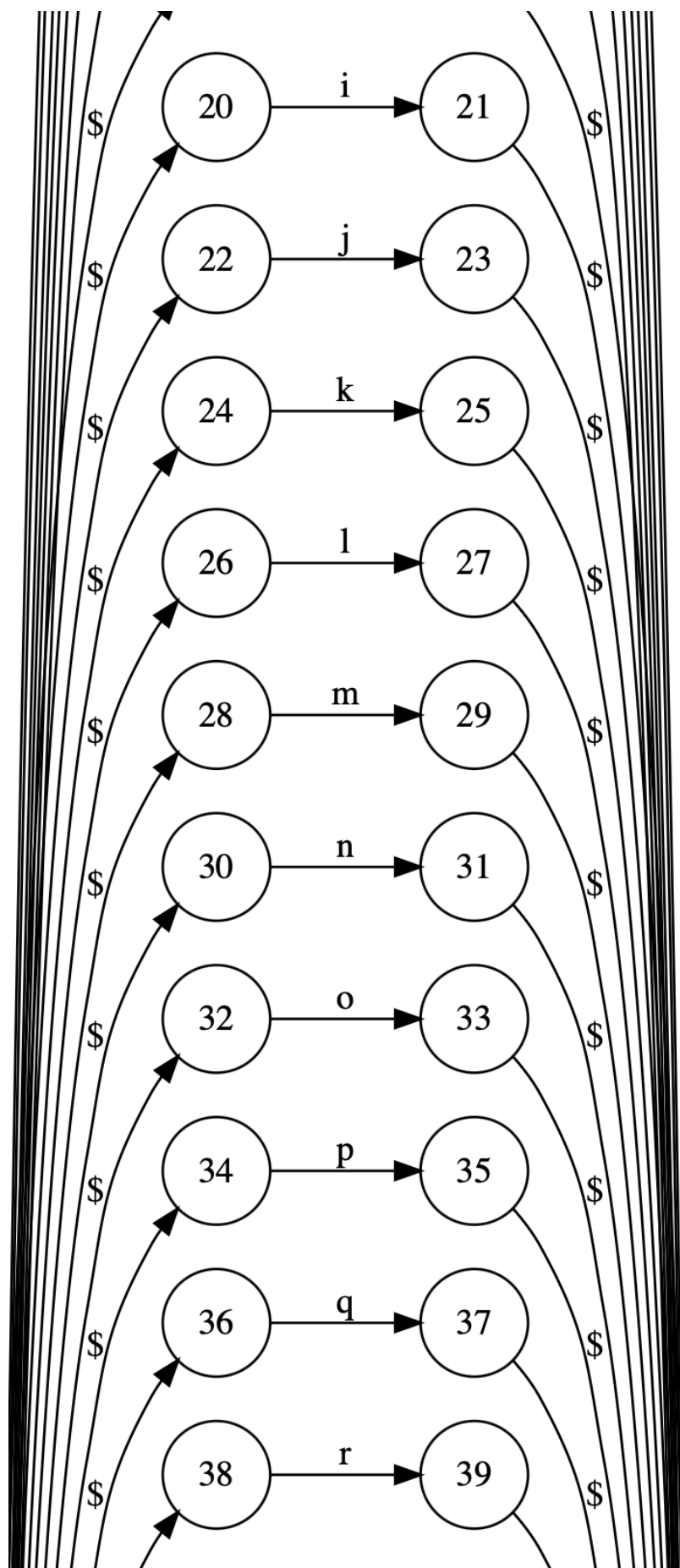
DFA for nonzero

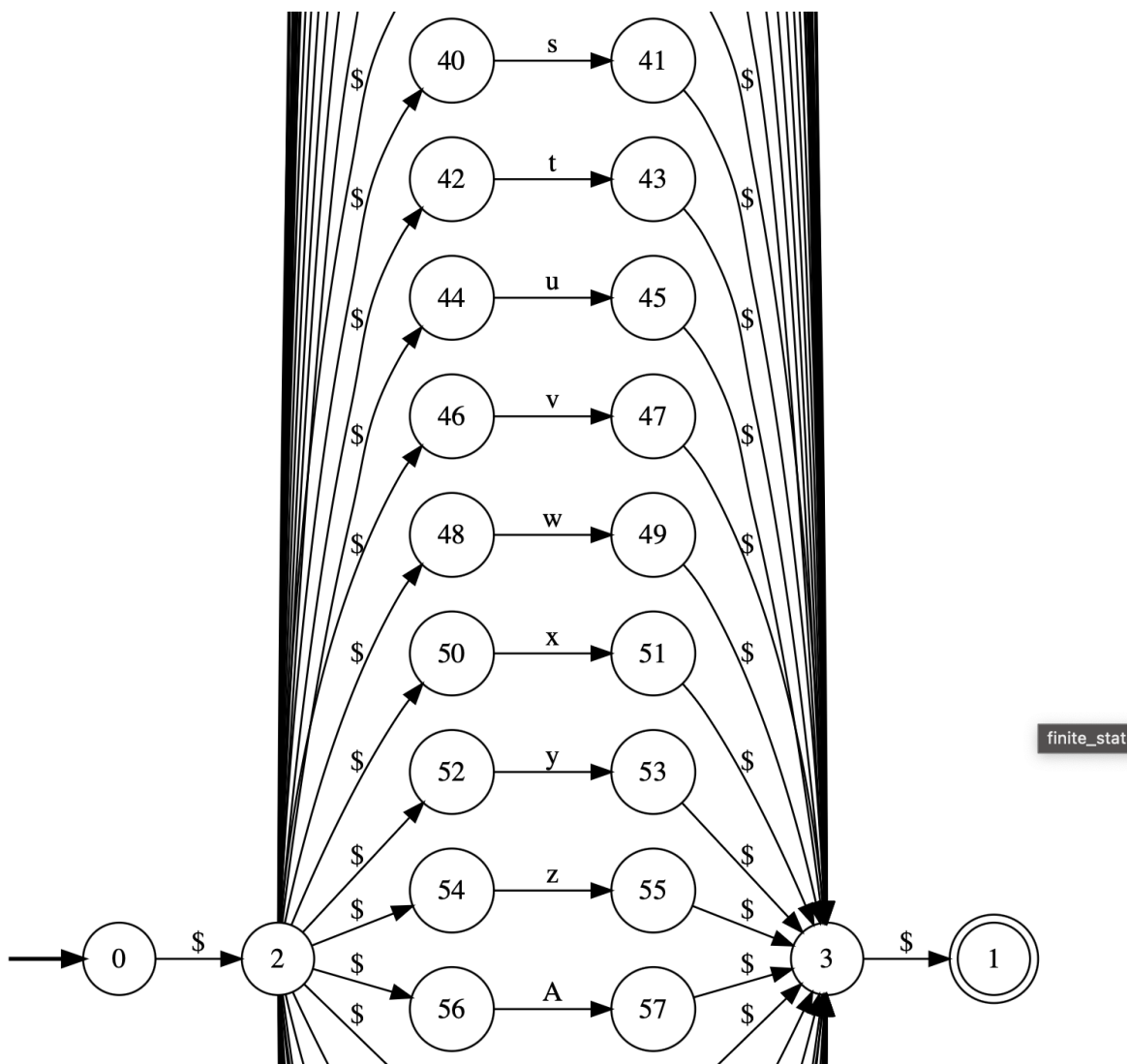


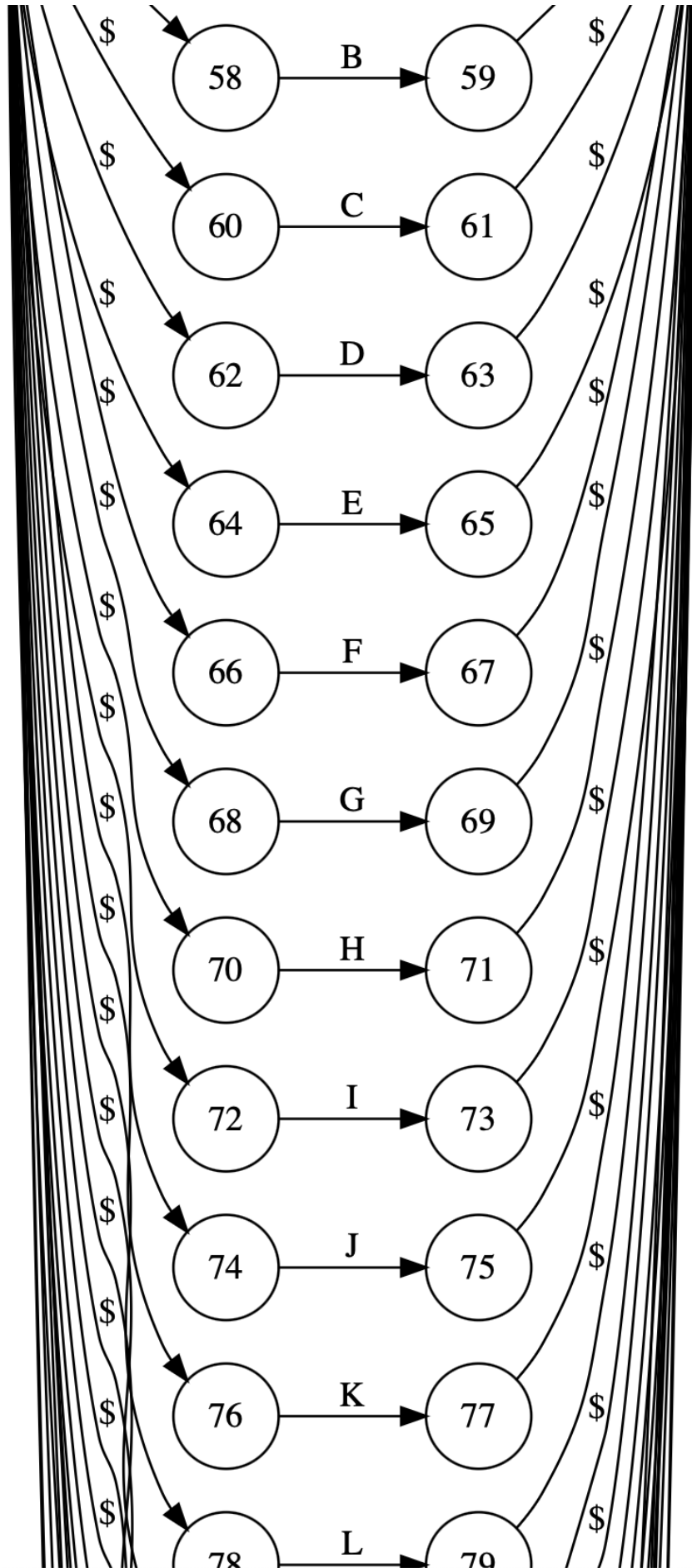
Regular expression for nonzero: $(1+2+3+4+5+6+7+8+9)$

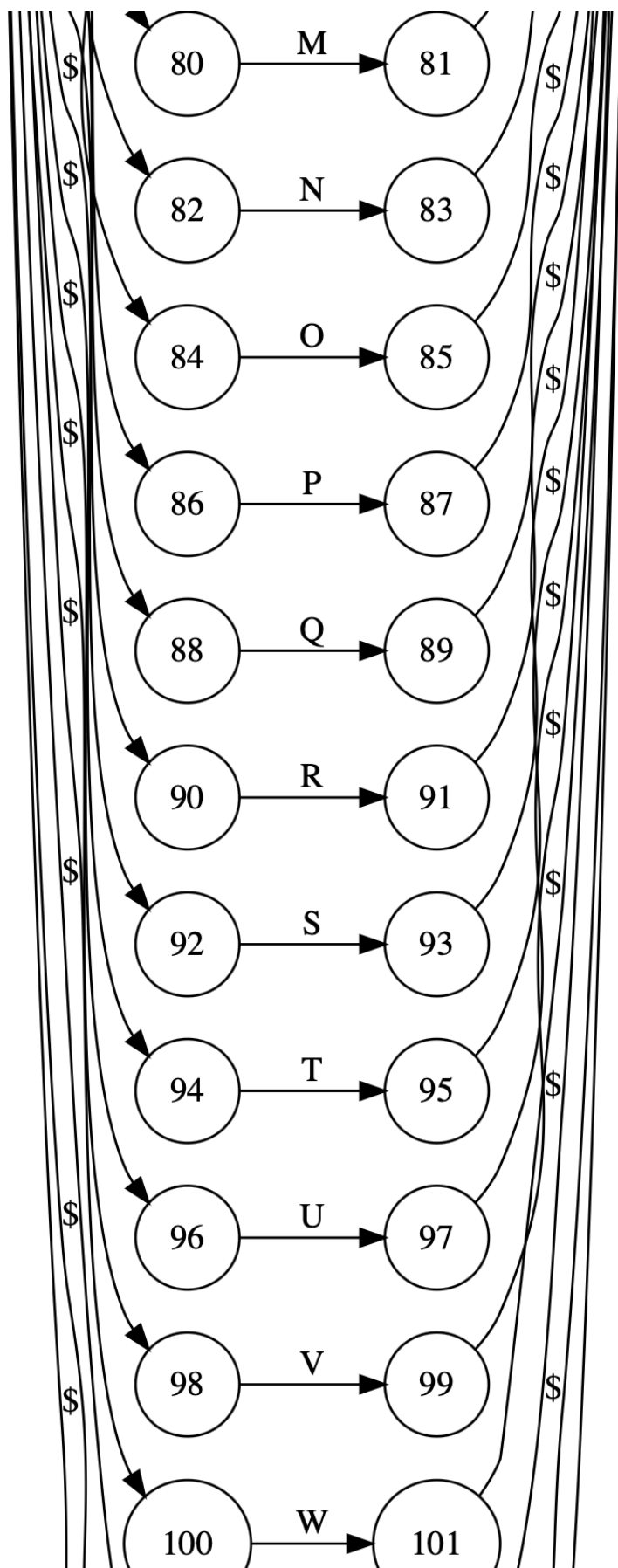
Finite state machine for letter

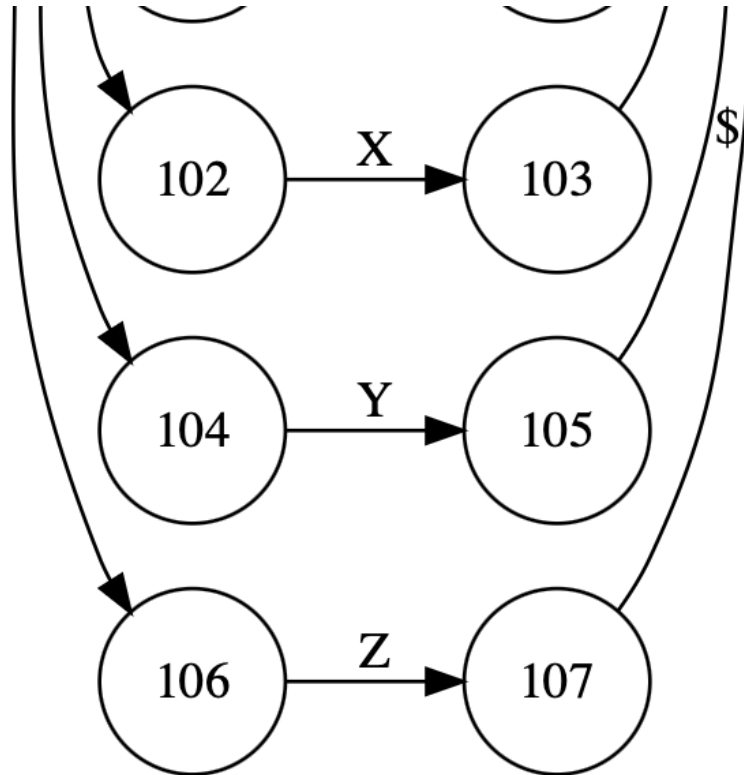












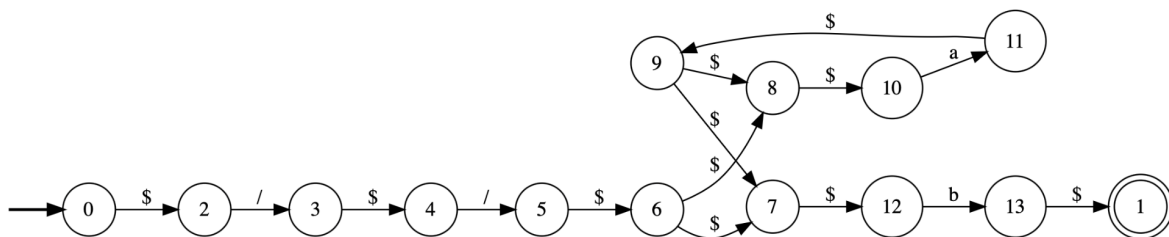
Regular expression for letter:

(a+b+c+d+e+f+g+h+i+j+k+l+m+n+o+p+q+r+s+t+u+v+w+x+y+z+A+B+C+D+E+F+G+H+I+J+K+L+M+N+O+P+Q+R+S+T+U+V+W+X+Y+Z)

Finite state machine for inline comment

A represents anything except a line break

B represents a line break

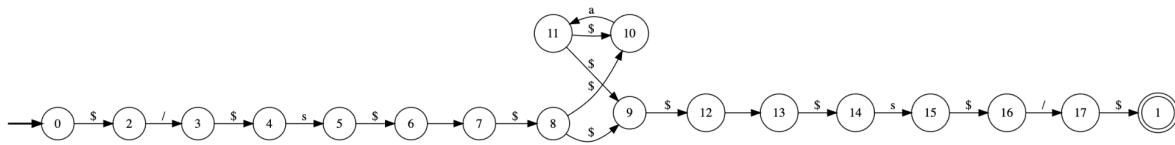


Regular expression for inline comment: `/(a)*b`

Finite state machine for block comment

S represents a *

A represents anything



Regular expression for block comment: `/s a* s/`

Section 3

I have a main file that is called `lexdriver` and is in charge of calling `getNextToken` and creating the `outlextokens/outlexerrors`. The actual logic for `getNextToken` is implemented in the class called `LexicalAnalyzer`. When an instance of that class is created, it will create a byte array with the bytes of the file as well as keep a current pointer. When `getNextToken` is called, it will get rid of all blankspaces until it can read a character. Based on what character it reads, it will go down a specific path that has been described in the dfa. For operations like `"="`, `">="` and `"=="` I used backtracking in order to make sure I don't skip over any characters. For longer tokens like `id`, `float` and `integer`, I read all the elements from the specific alphabet and I use `regex` afterwards to determine if it is a valid token (alternative 1).

I also have a `TokenType` enum which keeps track of all the types possible which makes it easier for building the token. Finally, I have a `Token` class that is responsible for creating the tokens and returning them to the user.

Section 4

I did my lexical analyzer in Java.

I used `StringBuilder` to construct the lexeme.

I used `Java.io` for file manipulation.

I used an enum for my token type (to be able to infer the type)

I used `Jest` for unit testing and to keep track of my code coverage.