Question 1. (20 points) Regular expressions and DFAs.

(a) (10 points) Give a regular expression that generates all strings with any combination of a's, b's, and c's such that the strings contain at least one a and at least one b. Except for this restriction, the letters in the string may appear in any order – in particular the required b may appear before or after the required a.

Fine print: You must restrict yourself to the basic regular expression operations covered in class and on homework assignments: rs, r|s,  $r^*$ , r+, r?, character classes like [a-cxy] and [^aeiou], abbreviations *name=regexp*, and parenthesized regular expressions. No additional operations that might be found in the "regexp" packages in various Unix programs, scanner generators like JFlex, or programming language libraries are allowed.

#### Any answer that generates the set of strings described in the problem is acceptable. Here is one:

```
[abc]*((a[ac]*b)|(b[bc]*a))[abc]*
```

(b) (10 points) Draw a DFA that accepts all valid strings of a's, b's, and c's that contain at least one a and at least one b (i.e., the set described above and generated by the regular expression from part (a)).



**Question 2.** (16 points) Ambiguity. Here is a small grammar for somewhat silly English sentences.

```
S ::= see the D T A \mid see the D A
D ::= orange \mid happy
T ::= duck \mid crab
A ::= walk \mid T walk
```

Is this grammar ambiguous? If so, give a proof that it is by showing two distinct parse trees or two different leftmost (or two different rightmost) derivations for some string generated by the grammar. If not, give an informal, but precise, argument why it is not ambiguous.

Notes: whitespace in the grammar is only for readability and is not part of the grammar or the strings generated by it. Each of the words like see or happy in a grammar production should be treated as a single terminal symbol, not as individual letters.

The grammar is ambiguous. Here are two distinct leftmost derivations of "see the orange duck walk"

S => see the D T A => see the orange T A => see the orange duck A => see the orange duck walk S => see the D A => see the orange A => see the orange T walk => see the orange duck walk

It also would be possible to show the ambiguity by showing the two distinct parse trees for this, or any other ambiguous, sentence.

**Question 3.** (30 points) The "they *still* can't think of *anything* original!!" LR parsing question. Here is a simple grammar for a language with a print and if statement and the terminal symbol x. The extra  $S' ::= S \$  rule needed to handle end-of-file in an LR parser has been added for you. As is usual, whitespace in the grammar is only for readability and is not part of the grammar or the strings generated by it. Terminal symbols like print or if are single symbols, not strings of letters.

0.	S' ::= S \$	(\$ is end-of-file)	2.	S ::= if E S
1.	S ::= print	Ε	3.	E ::= x

(a) (16 points) Draw the LR(0) state machine for this grammar. When you finish, you should number the states in the final diagram in whatever order you wish so you can use the state numbers in later parts of this question.



Question 3. (cont.) Grammar repeated from previous page for reference if needed:

0.	S' ::= S \$	(\$ is end-of-file)	2.	$S ::= \operatorname{if} E S$
1.	S ::= print	Ε	3.	E ::= x

(b) (10 points) Write the LR(0) parser tables for the LR parser in your answer to part (a).

	X	print	if	\$	E	S
0				acc		
1		s2	s4			<b>g0</b>
2	s5				g3	
3	r1	r1	r1	r1		
4	s5				<b>g6</b>	
5	r3	r3	r3	r3		
6		s2	s4			<b>g</b> 7
7	r2	<b>r2</b>	r2	r2		

(c) (4 points) Is this grammar LR(0)? Explain why or why not. Your answer should describe **all** of the problems that exist if the grammar is not LR(0) by identifying the relevant state number(s) in your answers to parts (a) and (b) and the specific issues in those state(s) (i.e., something like "state 47 has a shift-reduce conflict if the next input is  $f \circ \circ$ ", but with, of course, state numbers and correct details from your parser). If the grammar is LR(0) you should explain why (this can be brief).

Yes, it is LR(0). There are no shift/reduce or reduce/reduce conflicts in any states of the LR(0) state machine or parser table.

**Question 4.** (16 points) LL parsing. Here is a small grammar that generates strings of the letters a, m, n, and o. As usual, whitespace in the grammar rules is only for readability and not part of the generated strings.

- 1.  $R := aTC \circ | \circ$
- 2.  $T ::= m \mid \varepsilon$
- 3.  $C ::= n | \varepsilon$

(a) (8 points) Complete the following table showing the FIRST and FOLLOW sets and nullable for each of the non-terminals in this grammar:

	FIRST	FOLLOW	nullable
R	a, o		no
Т	m	n, o	yes
С	n	0	yes

Note: since (an end-of-file marker) is not included in the grammar, it really is not part of the set FOLLOW(*R*), but if it was included as part of that set in an answer there was no deduction.

(b) (8 points) Is this grammar, as written, suitable for constructing a top-down LL(1) predictive parser? If it is, your answer should give a technical explanation why it is. If not, your answer should give a technical explanation describing the problem or problems with this particular grammar that prevent it from being suitable for a LL(1) predictive parser.

Yes. For R, both productions start with different terminal symbols so we can always pick the right right-hand-side when expanding R. For T and C, their FIRST sets do not contain any symbols in their FOLLOW sets, so we can always choose between the production that expands the T or C to a terminal symbol or the epsilon production that erases the T or C from the derived string.

**Question 5.** (18 points) Semantics. Suppose we have the following assignment statement in a MiniJava program:

$$p = 2*x.f(a[i]);$$

(a) (8 points) Draw an abstract syntax tree (AST) for this statement at the bottom of this page. You should use appropriate names for AST nodes and have an appropriate level of abstraction and structural detail similar to the AST nodes in the MiniJava project AST classes, but don't worry about matching the exact names or details of classes or nodes found in the MiniJava code.

(b) (10 points) Annotate your AST by writing next to the appropriate nodes the checks or tests that should be done in the static semantics/type-checking phase of the compiler to ensure that this assignment statement does not contain errors. If a particular check or test applies to multiple nodes, you can write it once and indicate which nodes it applies to, as long as your meaning is clear and readable. You may assume that int is the only numeric type in the language, but remember that MiniJava also contains boolean and object (class) types.



Here are the semantics checks we should make:

- All variables are declared and visible in the current scope
- Assignment (=): verify p (left operand) designates an assignable location; verify the type of the expression (\*) is assignment-compatible with type of p
- Multiplication (\*): verify types of operands are int; type of \* is int
- Call: verify the type of the expression x has a one-argument function f whose parameter type is assignment-compatible with the type of the subscript ([]) expression (int here). The result type of the call node is the declared result type of f
- Subscript operator ([]): verify left operand has type array of int (the only possible array type in MiniJava). Verify the type of the right subscript operand is int. The type of the subscript [] node is int

Note: there are other possible representations for the AST. Answers that differed from the above solution received full credit if the AST choices were reasonable and at the right level of detail.