

# EPFL Computer Language Processing Exam (CS-320)

November 2022

Keep in mind the following:

1. **The exam is from 13:15 to 15:00. Do not open the exam until we tell you to.**
2. **Place your CAMIPRO card on your desk.**
3. **Put all electronic devices in a bag away from the bench.**
4. **Write your final answers using a permanent pen (no graphite, no “frixion” pen).**
5. **The only written material you are allowed to use in the exam is one sheet (two A4 pages) of any content that you prepare, either hand-written or printed.**
6. **Fill in all your answers on the given exam sheet. Do not submit additional sheets. Do not unstaple the sheets.**
7. Each question is scored independently. When you need to circle answers, then circling the correct set of answers gives you full points. If you circle some of the wrong answers or do not circle all the correct answers, you will obtain zero points on that question.
8. We advise you to first solve questions that you find easier.
9. The maximal number of points on the exam is **30**.

## How to circle answers:

- ① One correct solution; use hollow circle to circle the number only
2. Wrong answer; leave as is
- ③ Another correct solution; use hollow circle to circle the number only
4. Another wrong answer; leave as is

Write your full name on the line below:

---

Write your SCIPER on the line below:

---

## 1 Longest-Match Rule (2pt)

Consider the lexical analyzer accepting sequences over the alphabet  $A = \{0, 1, .\}$ , with token classes defined by following regular expressions (where float constants have only one decimal place):

**float** ::=  $(0 | 1)^* \cdot (0 | 1)^?$

**int** ::=  $1 (0 | 1)^*$

**zero** ::=  $0$

Note that in the notation above, the parentheses  $()$ ,  $()$  (as well as  $*$ ,  $?$ ,  $|$  are not part of the input alphabet but part of the notation for regular expressions.

Circle numbers next to correct statements regarding the rules above, assuming the longest match rule is applied:

1. **100.00** will be tokenized to sequence **float**
2. **100.00.1** will be tokenized to sequence **float, zero, float**
3. **10.10** will be tokenized to sequence **float, int**
4. **011.0** will be tokenized to sequence **zero, float**

[Solution: all statements are wrong.]

## 2 Epsilon Closure (4pt)

Consider the following representation of a NFA (non-deterministic finite automaton) in Scala

```
abstract class Transition
type State = Int
case class Epsilon(to: State) extends Transition
case class CharTransition(char: Char, to: State) extends Transition
case class NFA(initial: State, accepting: Set[State], delta: Map[State, Set[Transition]]) {
  def transitions(from: State): Set[Transition] = delta.getOrElse(from, Set.empty)

  def successors(from: State, visited: Set[State]): Set[State] =
    transitions(from).collect { case Epsilon(to) if !visited.contains(to) => to }

  def epsilonClosure(states: Set[State]) =
    def epsilonClosureRec(states: Set[State], acc: Set[State]): Set[State] = ???
    epsilonClosureRec(states, Set.empty)
}
```

The alphabet of our NFA is the Char type in Scala. Each state in the NFA is represented by an integer and a set of transitions to other states. A transition results in a destination state and is labeled either by  $\varepsilon$  or a character  $c : \text{Char}$ . We would like to complete the function `epsilonClosure` (which finds states reachable through epsilon transitions), so that it matches the definition seen in class. Circle all correct implementations of `epsilonClosureRec` among these, if any:

- ① **def** epsilonClosureRec(states: Set[State], acc: Set[State]): Set[State] =  
states.foldLeft(acc ++ states)((foldAcc, state) =>  
epsilonClosureRec(successors(state, foldAcc), foldAcc))
- ② **def** epsilonClosureRec(states: Set[State], acc: Set[State]): Set[State] =  
states.foldLeft(acc ++ states)((foldAcc, state) =>  
epsilonClosureRec(successors(state, acc), foldAcc))

Now circle all correct implementations among these two (if any):

- ① **def** epsilonClosureRec(states: Set[State], acc: Set[State]): Set[State] =  
**val** expansion = states.flatMap(state => successors(state, acc))  
**if** expansion.size > 0 **then** epsilonClosureRec(expansion, acc ++ states) **else** acc ++ states
1. **def** epsilonClosureRec(states: Set[State], acc: Set[State]): Set[State] =  
**val** expansion = states.flatMap(state => successors(state, acc))  
**if** expansion.size > 0 **then** epsilonClosureRec(expansion, acc ++ states) **else** acc

See next page for excerpts from the standard library reference for `foldLeft` and `collect`.

- **def** foldLeft[B](z: B)(op: (B, A) ⇒ B): B

Applies a binary operator to a start value and all elements of this set, going left to right.

- **B** the result type of the binary operator.
- **z** the start value.
- **op** the binary operator.
- **returns** the result of inserting **op** between consecutive elements of this set, going left to right with the start value **z** on the left:  $op(\dots op(z, x_1), x_2, \dots, x_n)$  where  $x_1, \dots, x_n$  are the elements of this set. Returns **z** if this set is empty.

- **def** collect[B](pf: PartialFunction[A, B]): Set[B]

Builds a new set by applying a partial function to all elements of this set on which the function is defined.

- **B** the element type of the returned set.
- **pf** the partial function which filters and maps the set. It acts both as a filter (by filtering out elements for which the function is not defined) and as a mapping.
- **returns** a new set resulting from applying the given partial function **pf** to each element on which it is defined and collecting the results. The order of the elements is preserved.

---

### LL(1) Parsing

---

First, consider the following grammar with non-terminals S,A and terminals **EOF**, ), (, ], [.

$S ::= A \text{ EOF}$

$A ::= ( A ) A \mid A [ A ] \mid \varepsilon$

### 3 Circle all true statements about the grammar above (2pt)

1. “[ ] ( [ ] )” is accepted by the grammar
2. The grammar is LL(1)
3. The grammar is ambiguous
4.  $\text{NULLABLE}(A) == \text{true}$
5.  $\text{NULLABLE}(S) == \text{true}$

### 4 Circle the correct answer (1pt)

1.  $\text{FIRST}(S) == \{ \text{EOF} \}$
2.  $\text{FIRST}(S) == \{ (, [ \}$
3.  $\text{FIRST}(S) == \{ (, ), \text{EOF} \}$
4.  $\text{FIRST}(S) == \{ (, [, \text{EOF} \}$
5.  $\text{FIRST}(S) == \{ (, ), [, ], \text{EOF} \}$

### 5 Circle the correct answer (1pt)

1.  $\text{FOLLOW}(A) == \{ ), ] \}$
2.  $\text{FOLLOW}(A) == \{ ), ], \text{EOF} \}$
3.  $\text{FOLLOW}(A) == \{ (, [, ), ] \}$
4.  $\text{FOLLOW}(A) == \{ (, [, ], \text{EOF} \}$
5.  $\text{FOLLOW}(A) == \{ (, [, ), ], \text{EOF} \}$

Answer: none of the above are correct, instead,  $\text{FOLLOW}(A) = \{ ), [, ], \text{EOF} \}$

## 6 Fill the LL(1) parsing table for the above grammar (2pt)

If referring directly to right-hand sides, the table looks as follows:

	<b>EOF</b>	(	)	[	]
S	A EOF	A EOF		A EOF	
A	$\epsilon$	(A) A, A [ A ]	$\epsilon$	A [ A ], $\epsilon$	$\epsilon$

If we number the right hand sides for each non-terminal in order:

$$\begin{aligned}
 S &::= A \mathbf{EOF} \\
 &\quad (1) \\
 A &::= (A)A \quad | \quad A[A] \quad | \quad \epsilon \\
 &\quad (1) \quad | \quad (2) \quad | \quad (3)
 \end{aligned}$$

then the table is the following one:

	<b>EOF</b>	(	)	[	]
S	1	1		1	
A	3	1,2	3	2,3	3

Now consider the following **different grammar**:

$$S ::= B \mathbf{EOF}$$

$$B ::= ( B ) B \mid \mathbf{empty}$$

where **empty** is a special keyword.

A student writes a parser in Scala for this language and implements the following case classes

```
case class B(children: Either[(parOpenToken, B, parCloseToken, B),
                             EmptyKW])
```

where parOpenToken, parCloseToken, EmptyKW are types representing the tokens (, ), **empty**.

## 7 Circle all the correct answers about the above class B (1pt)

1. The case class B *is suitable* as a node of an *abstract syntax tree (AST)* because it can be traversed recursively and is the most abstract representation sufficient to interpret and compile the program.
- ② The case class B *is suitable* as a node of a *parse tree* because the children of each node correspond to the right-hand side of the corresponding grammar rule.
3. The case class B *is not suitable* as a node of a parse tree because it is too abstract for this purpose.
- ④ The case class B *is not suitable* as a node of an *abstract syntax tree (AST)* because it contains unnecessary terminals.

---

Now consider the **following grammar**:

$S ::= C \text{ EOF}$

$C ::= \langle \text{IntegerLiteral} \rangle + C \mid \langle \text{IntegerLiteral} \rangle$

with terminals  $\langle \text{IntegerLiteral} \rangle$ ,  $+$ , **EOF**.

The student implements the following case class:

**case class** C(children:Either[(Int, C), Int])

## 8 Circle all correct answers about the above case class C (1pt)

- ① The class C is suitable as a node of an abstract syntax tree (AST) because it represents the information about the program that suffices to interpret and compile it.
2. The class C is suitable as a node of the parse tree because the children of each node correspond to the right-hand side of the corresponding grammar rule.
3. The class C is not suitable as a node of the AST because it contains unnecessary terminals.
- ④ The class C does not represent a parse tree node because it does not represent all terminals.

## 9 Properties of Grammars (1pt)

Circle all correct answers.

<p><math>S ::= ( S )</math></p> <p><math>R ::= ( T ) \mid T</math></p> <p><math>T ::= X \text{ or } X \mid X Z X \mid \varepsilon</math></p> <p><math>X ::= \text{true} \mid \text{false} \mid ( X )</math></p> <p><math>Y ::= \text{or}</math></p> <p><math>Z ::= \text{and}</math></p>	<ol style="list-style-type: none"> <li>1. There are exactly 3 unproductive symbols</li> <li>2. There is exactly 1 unreachable symbol</li> <li>3. There are exactly 3 unit productions</li> <li><input checked="" type="radio"/> 4. There is exactly 1 unit production</li> <li>5. There are exactly 2 unproductive symbols</li> <li>6. There are exactly 4 unreachable symbols</li> <li><input checked="" type="radio"/> 7. There are exactly 5 unreachable symbols</li> <li><input checked="" type="radio"/> 8. Z is unreachable</li> </ol>
--	--

Terminals are:  $)$ ,  $($ ,  $\text{or}$ ,  $\text{true}$ ,  $\text{false}$ ,  $\text{and}$ . Starting symbol is  $S$ .

## 10 Is this grammar in Chomsky Normal Form (CNF)? (1pt)

Circle all correct answers.

<p><math>S ::= \varepsilon \mid M N \mid O C</math></p> <p><math>S ::= \langle \text{IntegerLiteral} \rangle</math></p> <p><math>O ::= (</math></p> <p><math>C ::= )</math></p> <p><math>M ::= -</math></p> <p><math>N ::= \langle \text{IntegerLiteral} \rangle</math></p>	<ol style="list-style-type: none"> <li>1. No, because some rules do not have a valid number of symbols on the right-hand side</li> <li><input checked="" type="radio"/> 2. Yes, because only the starting symbol has <math>\varepsilon</math> on the right hand side and every other right hand side either has exactly two non-nonterminals, or exactly one terminal.</li> <li>3. No, because the starting symbol has an <math>\varepsilon</math></li> <li>4. No, because there is a unit production</li> <li>5. No, because integer literals are not allowed in CNF</li> <li>6. No, because the starting symbol should not have a terminal on its right-hand side</li> </ol>
---	--

Terminals are:  $)$ ,  $($ ,  $-$ ,  $\langle \text{IntegerLiteral} \rangle$ . Starting symbol is  $S$ .

## 11 Is this grammar in Chomsky Normal Form (CNF)? (1pt)

Circle all correct answers.

$S ::= \varepsilon \mid C$ $C ::= C_1 B$ $C_1 ::= A X$ $X ::= \text{true} \mid \text{false}$ $A ::= [$ $B ::= ]$	<ol style="list-style-type: none"> <li>1. No, because some rules do not have a valid number of symbols on the right-hand side</li> <li>2. Yes, because every rule has at most 2 symbols on the right-hand side, all terminals have an associated non-terminal, only the starting symbol is nullable, and there are no unit productions, unproductive or unreachable symbols</li> <li>③. No, because there is a unit production</li> <li>4. No, because not all terminals are on the right-hand side</li> <li>5. No, because the starting symbol has an <math>\varepsilon</math></li> </ol>
---	--

Terminals are: ], [, true, false. Starting symbol is S.

## 12 Fill the CYK parse table for the word “aba” and the following grammar. (2pt)

$S ::= A B \mid t_1 A \mid c$ $t_1 ::= a$ $A ::= t_1 A \mid c$ $B ::= t_2 B \mid b$ $t_2 ::= b$	<table border="1" style="margin: auto;"> <tr> <td style="padding: 5px;">a</td> <td style="padding: 5px;">b</td> <td style="padding: 5px;">a</td> </tr> <tr> <td style="padding: 5px;"><math>t_1</math></td> <td style="padding: 5px;"><math>B, t_2</math></td> <td style="padding: 5px;"><math>t_1</math></td> </tr> <tr> <td colspan="3" style="padding: 5px;"> </td> </tr> <tr> <td colspan="2" style="padding: 5px;"> </td> <td style="padding: 5px;"> </td> </tr> <tr> <td colspan="3" style="padding: 5px;"> </td> </tr> </table>	a	b	a	$t_1$	$B, t_2$	$t_1$									
a	b	a														
$t_1$	$B, t_2$	$t_1$														

Terminals are: a, b, c.

**13 What grammars could we obtain from the Initial Grammar after reducing the number of symbols on all right-hand sides to at most two? (2pt)**

Assume the algorithm that introduces no more non-terminals than the one presented in the lecture. Circle numbers for all grammars that could be obtained.

<p>Initial Grammar:</p> $S ::= (L)   a$ $L ::= L, S   S$	<p>1. Grammar:</p> $S ::= X   a$ $L ::= L V   S$ $V = , S$ $X ::= (L)$	<p>2. Grammar:</p> $S ::= O L C   a$ $L ::= L V S   S$ $O = (L$ $C = )$ $V = ,$
<p>③ Grammar:</p> $S ::= O C   a$ $L ::= V S   S$ $O = (L$ $C = )$ $V = L ,$	<p>4. Grammar:</p> $S ::= Y   a$ $L ::= J   S$ $X = (L$ $Y = X)$ $Z = L ,$ $J = Z S$	<p>5. Grammar:</p> $S ::= X   a$ $L ::= V S   S$ $V = L ,$ $X = O C$ $O = ($ $C = L K$ $K = )$

Terminals are: a, ), ( and the comma , .

Solution note: answer set {3} was accepted, and, to be generous, answer set {3, 4, 5} as well.

## 14 Typing Rules (3pt)

See next page.

## 15 Type Inference (3pt)

1.  $x \Rightarrow y \Rightarrow y(z \Rightarrow 6) + y(7)$
- ②.  $g \Rightarrow f \Rightarrow x \Rightarrow g(f(x))$
- ③.  $x \Rightarrow y \Rightarrow ((z \Rightarrow y), y)$
- ④.  $g \Rightarrow f \Rightarrow x \Rightarrow g(f(x)) + f(g(x)) + x$

## 16 Unification Algorithm (3pt)

1. 

$\tau = (\tau_{10}, \tau_1)$
$\tau = (\tau_2, \tau_{20})$
$\tau_3 = (\tau_1, \tau_2)$
2. 

Elimination of $\tau$ : $\tau = (\tau_{10}, \tau_1)$ $(\tau_{10}, \tau_1) = (\tau_2, \tau_{20})$ $\tau_3 = (\tau_1, \tau_2)$
--
3. 

Decomposition of pairs: $\tau = (\tau_{10}, \tau_1)$ $\tau_{10} = \tau_2, \tau_1 = \tau_{20}$ $\tau_3 = (\tau_1, \tau_2)$
---
4. 

Elimination of $\tau_{10}$ : $\tau = (\tau_2, \tau_1)$ $\tau_{10} = \tau_2, \tau_1 = \tau_{20}$ $\tau_3 = (\tau_1, \tau_2)$
---

(argument)       $t$  : 

$(\tau_2, \tau_1)$
--------------------

(result)       $\text{swap}(t)$  : 

$(\tau_1, \tau_2)$
--------------------

$$\frac{\frac{\frac{(x, Int) \in \Gamma}{\Gamma \vdash x : Int} \quad \frac{(u, Int) \in \Gamma}{\Gamma \vdash u : Int}}{\Gamma \vdash x < u : Bool} \quad \frac{\frac{\frac{(v, Int) \in \Gamma}{\Gamma \vdash v : Int} \quad \frac{(x, Int) \in \Gamma}{\Gamma \vdash x : Int}}{\Gamma \vdash v < x : Bool} \quad \frac{(v, Int) \in \Gamma}{\Gamma \vdash v : Int} \quad \frac{(x, Int) \in \Gamma}{\Gamma \vdash x : Int}}{\Gamma \vdash (v < x) \vee (v < x) \text{ else } x : Int}}{\Gamma \vdash \mathbf{if}(x < u)u \text{ else } \mathbf{if}(v < x)v \text{ else } x : Int}$$