

# Informatics 2A: Tutorial Sheet 4 - SOLUTIONS

MARY CRYAN

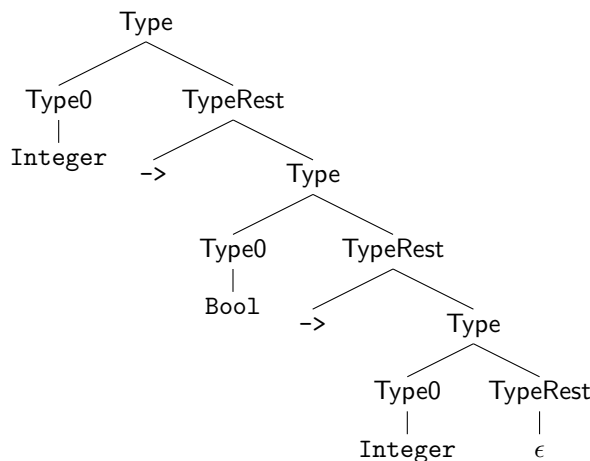
1. A possible LL(1) grammar is as below (note the distinction between  $\epsilon$  and  $\epsilon!$ ). It's wise to also talk to them about the similarity between this and what we did to get an LL(1) grammar for the arithmetic expressions:

$$\begin{aligned}
 \text{RegExp} &\rightarrow \text{RegExp1 PlusOps} \\
 \text{PlusOps} &\rightarrow \epsilon \mid + \text{RegExp1 PlusOps} \\
 \text{RegExp1} &\rightarrow \text{RegExp2 ConcatOps} \\
 \text{ConcatOps} &\rightarrow \epsilon \mid \text{RegExp2 ConcatOps} \\
 \text{RegExp2} &\rightarrow \text{RegExp3 StarOps} \\
 \text{StarOps} &\rightarrow \epsilon \mid * \text{StarOps} \\
 \text{RegExp3} &\rightarrow \text{Atom} \mid ( \text{RegExp} ) \\
 \text{Atom} &\rightarrow \text{sym} \mid \emptyset \mid \epsilon
 \end{aligned}$$

For the sake of completeness, the parse table for this is:

	sym/ $\emptyset$ / $\epsilon$	+	*	(	)	\$
RegExp	RegExp1 PlusOps			RegExp1 PlusOps		
PlusOps		+RegExp1 PlusOps				$\epsilon$ $\epsilon$
RegExp1	RegExp2 ConcatOps			RegExp2 ConcatOps		
ConcatOps	RegExp2 ConcatOps	$\epsilon$		RegExp2 ConcatOps		$\epsilon$ $\epsilon$
RegExp2	RegExp3 StarOps			RegExp3 StarOps		
StarOps	$\epsilon$	$\epsilon$	* StarOps	$\epsilon$		$\epsilon$ $\epsilon$
RegExp3	Atom			(RegExp)		
Atom	sym/ $\emptyset$ / $\epsilon$					

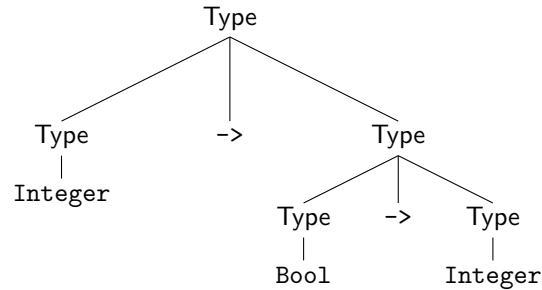
2. (a) The parse tree is:



- (b) By *abstract syntax tree* we mean the intended parse tree of the abstract syntax grammar. In this case it is the parse tree that implicitly brackets the type as

$$\text{Integer} \rightarrow (\text{Bool} \rightarrow \text{Integer})$$

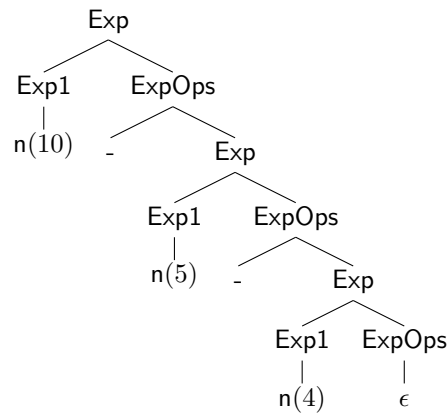
which is the correct bracketing by the Haskell convention that the  $\rightarrow$  operation associates to the right. The required tree is:



(c)

$$\begin{aligned} \text{Exp} &\rightarrow \text{Exp1 ExpOps} \\ \text{Exp1} &\rightarrow n \mid (\text{Exp}) \\ \text{ExpOps} &\rightarrow \epsilon \mid + \text{Exp} \mid - \text{Exp} \end{aligned}$$

(d) The parse tree is:

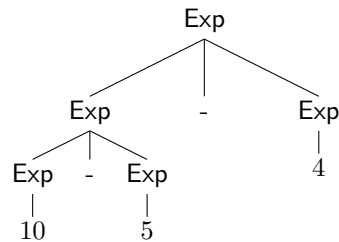


(At the numeral leaves, the labels on the parse tree are strictly speaking just  $n$ , but it's informative to display the actual lexical tokens in question.)

(e) This time, standard conventions require the abstract syntax tree that implicitly brackets the expression as

$$(10 - 5) - 4$$

This leads to the abstract syntax tree below.



(f) The difference is that structurally similar concrete parse trees are dealt with in different ways in the translation to abstract syntax tree. The concrete syntax does not distinguish between right-associativity and left-associativity of binary operations. Rather, irrespective of the associativity properties of the operations, expressions are concretely parsed as a *list* of top-level operations. The right-associativity of  $\rightarrow$  is catered for by translating concrete parse trees to abstract syntax trees in one way. The left-associativity of  $+$  and  $-$  is handled by translating in a different way.

The main point to take away from this is: when using LL(1) parse technology, the choice of whether to implement infix operations as left-associative or right-associative is made in the translation from concrete syntax to abstract syntax, it is not made in the formulation of the grammar for the concrete syntax.

3. This question is really to get the students to think about what they have learned in the “Fixing Grammars” lecture. The actual “work” they need to do to answer this is very minor.

They will remember the algorithm of Lecture 13 which takes a CFG and converts it to Chomsky Normal Form. This has 4 steps: (i) removing  $\epsilon$ -productions, (ii) removing unit-productions  $X \rightarrow Y$  where  $Y$  is non-terminal, (iii) removing terminal symbols from “mixed” right-hand sides of production rules (by adding dummy non-terminals) and (iv) breaking up long (at least three) right-hand sides by introducing more dummy non-terminals.

We only need to carry out the first two of these steps in order to get an equivalent CFG for  $\mathcal{G}$  which is cycle-free (see those steps on slide 14 of lecture 13 if you want to discuss more about it). Also, if the original grammar generated the empty string, we add an extra non-terminal  $\widehat{S}$  which becomes the new start non-terminal, and two extra production rules

$$\widehat{S} \rightarrow \epsilon \mid S.$$

Note that after having carried out step (i) of CNF, we are assured that every right-hand side of every terminal is non-empty, and for each non-terminal  $X$  in  $N$ , there is no derivation  $X \Rightarrow^* \epsilon$  (we make a note of whether  $S$  could derive the empty string, and will correct for this at the end). However, after (i) we may still have productions which have a single non-terminal on the right-hand side, and these may give rise to “cycles” as defined in the question.

Step (ii) of CNF conversion eliminates all unit production rules with a non-terminal on the right-hand side, resulting in a grammar where every production rule is either  $X \rightarrow a$  for a single non-terminal, or alternatively is  $X \rightarrow x_1 \dots x_k$  for  $k \geq 2$ . Since the non-terminals no-longer can generate an empty string, this “at least two symbols” constraint on the right-hand-side means we can never generate  $X \Rightarrow^+ X$  by introducing initial non-terminals and then “vanishing” them.

The adjustment given above for the case where  $\epsilon \in L(\mathcal{G})$  does not induce a cycle (since we use a dummy start non-terminal for this).