# 1 2004 Paper 4 Question 1

(a) A context-free grammar can be formally defined as a 4-tuple. Give a precise statement of what the components are

$$G = (N, T, P, S)$$

- $G$ is the grammar

- $N$ is the set of nonterminals

  A Nonterminal is an internal symbol. These represent concepts such as expressions or statements.

- $T$ is the set of terminals

  A Terminal is a token passed to the parser by the lexer. These may correspond to an individual literal or a sequence of literals. Terminals are indivisible. The input to any PDA is a sequence of terminals.

- $P \subseteq N \times (N \cup T)*$ is the set of productions

  A production is of the form $A \rightarrow \alpha$ and says that it is legal for any occurrence of $A$ to be replaced with $\alpha$ at any point.

- $S \in N$ is the start symbol

Note also that $N$, $T$ and $P$ are finite; and that $N \cup T = \emptyset$.

(b) Explain the difference between a grammar and the language it generates.

A grammar is a set of rules which is used to generate a language.

The language generated by a grammar is a set of strings.

Each grammar generates exactly one language, however a given language may be generated by many grammars.

A grammar is finite, while a language is infinite and a language is flat while a grammar is structured.
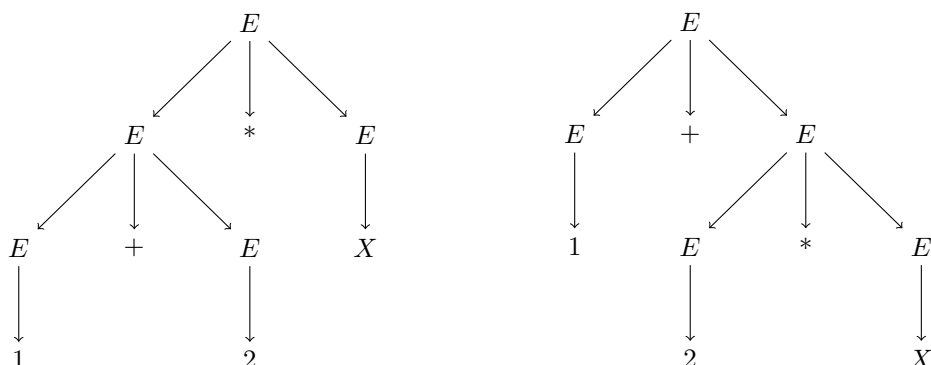
(c) Explain what makes a grammar ambiguous, with reference to the grammar which may commonly be expressed as a "rule"

$$E ::= 1 \mid 2 \mid X \mid E + E \mid E * E \mid -E$$

where $X$ is an identifier

A grammar is ambiguous if there exists any string for which there are multiple leftmost derivations for the grammar to generate that string. Consider the string $1 + 2 * X$ with the grammar above.

Under the grammar above, there are two possible parse trees for $1 + 2 * X$ and therefore the grammar is ambiguous.

(d) For the "rule" in part (c), give a formal grammar containing this "rule" and adhering to your definition in part (a).

$$G = (\{E\},$$
$$\{1, 2, X, *, +\},$$
$$\{(E, 1), (E', 2), (E', X), (E, E + E), (E, E * E), (E, -E)\},$$
$$E)$$

$$E ::= T\ E'$$
$$E' ::= +T\ E'\ |\ *T\ E'\ |\ \varepsilon$$
$$T ::= N\ |\ -T$$
$$N ::= 1\ |\ 2\ |\ X$$

(e) Give non-ambiguous grammars each generating the same language as your grammar in part (d) for the cases:

  (i) "$-$" is most tightly binding and "$+$" and "$*$" have equal binding power and associate to the left.

$$G_1 = (\{E, E', N, T\},$$
$$\{1, 2, X, *, +\},$$
$$\{(E, E'T), (E', E'T+), (E', E'T*), (E', \varepsilon), (T, N), (T, -T), (N, 1), (N, 2), (N, X)\},$$
$$E)$$

$$E ::= E'\ T$$
$$E' ::= E'\ T\ +\ |\ E'\ T\ *\ |\ \varepsilon$$
$$T ::= N\ |\ -T$$
$$N ::= 1\ |\ 2\ |\ X$$

  (ii) "$-$" is most tightly binding and "$+$" and "$*$" have equal binding power and associate to the right.

$$G_2 = (\{E, E', N, T\},$$
$$\{1, 2, X, *, +\},$$
$$\{(E, TE'), (E', +TE'), (E', *TE'), (E', \varepsilon), (T, N), (T, -T), (N, 1), (N, 2), (N, X)\},$$
$$E)$$

$$E ::= T\ E'$$
$$E' ::= +T\ E'\ |\ *T\ E'\ |\ \varepsilon$$
$$T ::= N\ |\ -T$$
$$N ::= 1\ |\ 2\ |\ X$$

  (iii) "$-$" binds more tightly than "$+$", but less tightly than "$*$", with "$+$" left-associative and "$*$" right-associative so that "$-a + -b * c * c + d$" is associated as "$((-a) + (-(b * (c * d)))) + d$".

$$G_3 = (\{E, E', T, F, F', N\},$$
$$\{1, 2, X, *, +\},$$
$$\{(E, E'T), (E', E'T+), (E', \varepsilon), (T, F), (T, -T), (F, NF'), (F', *NF'), (F', \varepsilon), (N, 1), (N, 2), (N, X)\},$$
$$E)$$

$$E ::= E' \; T$$
$$E' ::= E' \; T + \; | \; \varepsilon$$
$$T ::= F \; | \; - T$$
$$F ::= N \; F'$$
$$F' ::= *N \; F' \; | \; \varepsilon$$
$$N ::= 1 \; | \; 2 \; | \; X$$

(f) Give a simple recursive descent parser for your grammar in part (e)(iii) above which yields a value of type `ParseTree`. You may assume operations *mkplus*, *mktimes*, *mkneg* acting on type `ParseTree`.

```
type n = E | E' | T | F | F' | N

type t = Plus | Minus | Times | 1 | 2 | X | Epsilon

let parse ts =
        let parseE ts =
                let pt, ts = parseT ts in
                parseE' pt ts
        in
        let parseE' pt1 = function
                | Plus::ts -> (
                        let pt2, ts = parseT ts in
                        match ts with
                        | Plus::ts -> (
                                let pt2, ts = (parseT pt (Plus::ts)) in
                                parseE' (mkplus pt1 pt2) ts
                                )
                        | _ -> pt, ts
                        )
                | ts -> pt1
        in
        let parseT = function
                | Minus::ts -> (
                        let pt, ts = parseT ts in
                        (mkminus pt), ts
                        )
                | ts -> parseF ts
        in
        let parseF ts =
                let pt, ts = parseN ts in
                parseF' pt ts
        in
        let parseF' pt1 = function
                | Times::ts -> (
                        let pt2, ts = parseN ts in
                        parseF' (mktimes pt1 pt2) ts
                        )
                | ts -> pt1, ts
        in
        match parseE ts with
        | _, [] -> raise ParseException
        | pt, _ -> pt
```

## 2  2002 Paper 4 Question 2

The specification for a pocket-calculator-style programming language is as follows:

- Valid inputs consist either of an Expression followed by the ⟨enter⟩ button of of an Expression followed by ⟨store⟩ Identifier ⟨enter⟩;

- Expressions consist of Numbers and Identifiers connected with the binary operators ⟨+⟩, ⟨×⟩ and ⟨↑⟩ (in increasing binding power), with the Unary operators ⟨−⟩ and ⟨abs⟩, and possibly grouped with parentheses. Unary operators bind more strongly than ⟨+⟩ but weaker than ⟨×⟩ so that $-a + b$ means $(-a) + b$ but $-a \times b$ means $-(a \times b)$.

- Numbers consist of a sequence of at least one digit, possibly interspersed with exactly one decimal point, and possibly followed by an exponential marker "$e$" followed by a signed integer, e.g. $6.023e + 22$. Identifiers are sequences of lower-case letters.

(a) Give a Context-Free Grammar for the set of valid input sequences using names beginning with an upper-case letter for non-terminals. It should be complete in that you should go as far as to define e.g.

$$\textbf{Letter} ::= \textbf{a} \mid \textbf{b} \mid \textbf{c} \mid \ldots \mid \textbf{z}$$

$$
\begin{aligned}
\textbf{Start} &::= \textbf{Expression} \; \boxed{\text{enter}} \mid \textbf{Expression} \; \boxed{\text{store}} \; \textbf{Identifier} \; \boxed{\text{enter}} \\
\textbf{Expression} &::= \textbf{Unary OptExpression} \\
\textbf{OptExpression} &::= \boxed{+} \; \textbf{Unary OptExpression} \mid \varepsilon \\
\textbf{Unary} &::= \textbf{Times} \mid \boxed{-} \; \textbf{Unary} \mid \boxed{\text{abs}} \; \textbf{Unary} \\
\textbf{Times} &::= \textbf{Power OptTimes} \\
\textbf{OptTimes} &::= \boxed{\times} \; \textbf{Unary OptTimes} \mid \varepsilon \\
\textbf{Power} &::= \textbf{Value OptPower} \\
\textbf{OptPower} &::= \boxed{\uparrow} \; \textbf{OptUnary} \mid \varepsilon \\
\textbf{OptUnary} &::= \boxed{-} \; \textbf{OptUnary} \mid \boxed{\text{abs}} \; \textbf{OptUnary} \mid \textbf{Power} \\
\textbf{Value} &::= (\textbf{Expression}) \mid \textbf{Identifier} \mid \textbf{Number} \\
\textbf{Identifier} &::= \textbf{Letter OptIdentifier} \\
\textbf{OptIdentifier} &::= \textbf{Letter OptIdentifier} \mid \varepsilon \\
\textbf{Letter} &::= \textbf{a} \mid \textbf{b} \mid \textbf{c} \mid \ldots \mid \textbf{z} \\
\textbf{Number} &::= \textbf{Int OptInt OptDecimal OptSuffix} \mid \boxed{.} \; \textbf{Int OptInt OptSuffix} \\
\textbf{Int} &::= 0 \mid 1 \mid \ldots \mid 9 \\
\textbf{OptInt} &::= \textbf{Int OptInt} \mid \varepsilon \\
\textbf{OptDecimal} &::= \boxed{.} \; \textbf{OptInt} \mid \varepsilon \\
\textbf{OptSuffix} &::= \textbf{e Sign Int OptInt} \\
\textbf{Sign} &::= \boxed{+} \mid \boxed{-}
\end{aligned}
$$

(b) Indicate, giving brief reasoning, which non-terminals are appropriate to be processed using lexical analysis and for which using syntax analysis is proper.

It's appropriate to process **Value**, **Identifier**, **OptIdentifier**, **Letter**, **Number**, **Int**, **OptInt**, **OptDecimal**, **OptSuffix** and **Sign** in lexical analysis. This is because the langauge which these non-terminals can match is regular and there is no binding tightness to consider. Therefore, it's appropriate to process them during lexing.

(c) Give yacc or CUP input describing those elements deemed in part (b) to be suitable for syntax analysis. You need not give "semantic actions".

```
%token  Start  Expression  OptExpression  Unary  Times  OptTimes  Arrow  OptArrow

%right  '+'
%noassoc  '−'
%left  '*'
%left  '↑'

%%

E          :  v
           |  E + E
           |  − E
           |  E * E
           |  E ↑ E
           |  E
```