# Week 2: Context Free Grammars

* 1. Consider the following CFG $G$ with start symbol $R$:

$$
\begin{aligned}
R &\longrightarrow XRX \mid S \\
S &\longrightarrow aTb \mid bTa \\
T &\longrightarrow XTX \mid X \mid \epsilon \\
X &\longrightarrow a \mid b
\end{aligned}
$$

(a) What are the non-terminals of $G$?

(b) What are the terminals of $G$?

(c) Give three strings in $L(G)$.

(d) Give three strings not in $L(G)$.

(e) True or false: $T \rightarrow aba$.

(f) True or false: $T \rightarrow^* aba$.

(g) True or false: $T \rightarrow T$.

(h) True or false: $T \rightarrow^* T$.

(i) True or false: $XXX \rightarrow^* aba$.

(j) True or false: $X \rightarrow^* aba$.

(k) True or false: $T \rightarrow^* XX$.

(l) True or false: $T \rightarrow^* XXX$.

(m) True or false: $S \rightarrow^* \epsilon$.

* 2. Figure 1 contains a grammar for the Brischeme language from this week's lab sheet. Which of the following are valid Brischeme programs (i.e. strings in the language of that grammar)? You do not have to give the derivations (but you should work through them in your head).

(a) (define $x$ 32) (+ 4 $x$)

The syntax of the *Brischeme* programming language (ignoring whitespace) is given by the CFG with:

- Terminals: $0, 1, \ldots, 9$, $a, b, \ldots, z$, $A, B, \ldots, Z$, _, !, ?, $<$, $=$, $+$, $-$, $*$, $/$, (, ), define, lambda, not, or, and, #t, #f.

- Nonterminals: *Prog*, *Form*, *SExpr*, *Ident*, *Literal*, *Num*, *Bool*, *Digit*, *LChar*, *IdChar*, *Primop*.

- The production rules are:

$$
\begin{aligned}
\textit{Prog} \quad &::= \quad \textit{Form}^* \\[4pt]
\textit{Form} \quad &::= \quad \textit{SExpr} \\
&\quad\; | \quad (\text{ define } \textit{Ident } \textit{SExpr } ) \\[4pt]
\textit{SExpr} \quad &::= \quad \textit{Literal} \\
&\quad\; | \quad \textit{Ident} \\
&\quad\; | \quad (\textit{ SExpr SExpr}^* ) \\
&\quad\; | \quad (\textit{ Primop SExpr}^* ) \\
&\quad\; | \quad (\text{ lambda } (\textit{ Ident}^* ) \textit{ SExpr } ) \\[4pt]
\textit{Ident} \quad &::= \quad \textit{LChar IdChar}^* \\
\textit{LChar} \quad &::= \quad \text{a} \,|\, \text{b} \,|\, \cdots \,|\, \text{z} \\
\textit{IdChar} \quad &::= \quad a \,|\, b \,|\, \cdots \,|\, z \,|\, A \,|\, B \,|\, \cdots \,|\, Z \,|\, ! \,|\, ? \,|\, \_ \\[4pt]
\textit{Literal} \quad &::= \quad \textit{Bool} \,|\, \textit{Num} \\
\textit{Bool} \quad &::= \quad \text{\#t} \,|\, \text{\#f} \\
\textit{Num} \quad &::= \quad \textit{Digit Digit}^* \\
\textit{Digit} \quad &::= \quad 0 \,|\, 1 \,|\, \cdots \,|\, 9 \\[4pt]
\textit{Primop} \quad &::= \quad + \,|\, - \,|\, {}^* \,|\, / \,|\, \text{and} \,|\, \text{or} \,|\, \text{not} \,|\, < \,|\, =
\end{aligned}
$$

Figure 1: Brischeme (ignoring whitespace).

(b) (define *x* 32) (+ 4 *y*)

(c) (define *x* 32) (+ 4 5*x*)

(d) (define aC3! 32) (+ 4 aC3!)

(e) (define AC3! 32) (+ 4 AC3!)

(f) (+ (define *x* 32) *x* 4)

(g) (lambda *x* (+ *x* *x*))

(h) (define *f* (lambda (*x y*) (+ *x* (* 2 *y*))))

(i) (lambda (*x b*) (+ *x* (not *b*)))

(j) (lambda (*x b*) (if *x*))

(k) (lambda (*x b*) ((if *b* + -) 2 3))

(l) ((lambda () 3))

** 3. Design CFGs for the following programming language lexemes over the ASCII alphabet. You will find it convenient to use abbreviations like $\cdots$ to help present the expressions compactly.

(a) A *C program identifier* is any string of length at least 1 containing only letters ('a'–'z', lower and uppercase), digits ('0'–'9') and the underscore, and which begins with a letter or the underscore.

(b) An *integer literal* is any string taking one of the following forms:
  - a non-empty sequence of digits (decimal)
  - a non-empty sequence of characters from '0'–'9','a'–'e' (upper or lowercase) that are preceded by "0x" (hexadecimal)
  - a non-empty sequence of bits '0' and '1' that are preceded by "0b" (binary)

** 4. Give a CFG that describes the language of all subsequences of all permutations (i.e. no repetition) of the following keywords:

final static synchronized

** 5. Give CFGs for the following languages. The later parts are more difficult than 2-star.

(a) All odd length strings over $\{a, b\}$.

(b) All strings over $\{a, b\}$ that contain $aab$ as a substring.

(c) The set of strings over $\{a, b\}$ with more $a$ than $b$. Hint: every string $w$ with at least as many

3

*a* as *b* (possibly the same number of *a* as *b*) can be characterised inductively as follows. Either:

- *w* is just *a*
- or, *w* is of shape *avb* with *v* containing at least as many *a* as *b*
- or, *w* is of shape *bva* with *v* containing at least as many *a* as *b*
- or, *w* is of shape $v_1 v_2$ with $v_1$ and $v_2$ each separately containing at least as many *a* as *b*
- or, *w* is the empty string

(d) The complement of the language $\{a^n b^n \mid n \geq 0\}$ over $\{a, b\}$. Hint: express "not of shape $a^n b^n$ for some *n*" into one or more positive (i.e. without using *not* or similar) conditions.

(e) $\{v \# w \mid v, w \in \{a, b\}^* \text{ and the reverse of } v \text{ is a substring of } w\}$, over $\{a, b, \#\}$.

*** 6. An $\epsilon$-production is a rule of shape $X \longrightarrow \epsilon$ (for some nonterminal $X$). A unit production is a rule of shape $X \longrightarrow Y$ (for some non-terminals $X$ and $Y$). Consider the following grammar $G$ with start symbol $S$:

$$S \longrightarrow aSbb \mid T$$
$$T \longrightarrow bTaa \mid S \mid \epsilon$$

This grammar has two unit productions $S \longrightarrow T$ and $T \longrightarrow S$, and it has an epsilon production $T \longrightarrow \epsilon$.

Give a grammar with *no* unit productions and *no* $\epsilon$-productions for the language $L(G) \setminus \{\epsilon\}$.

*** 7. Give an English description of the language of the grammar in Q1.