

Final #1

Mark all correct answers in each of the following questions.

Unless stated otherwise, $G = (N, T, R, S)$ is a context-free grammar without useless letters.

4. (a) If G is ambiguous, then there may exist words $w \in T^*$ produced by a unique parse tree.
- (b) If G_1, G_2 are unambiguous grammars and $L(G_1) \cap L(G_2) = \emptyset$, then it is possible to construct an unambiguous grammar G such that $L(G) = L(G_1) \cup L(G_2)$.
- (c) Suppose $G_i = (N_i, T, R_i, S_i), i = 1, 2$, where $N_1 \cap N_2 = \emptyset$. Let $G = (N_1 \cup N_2 \cup \{S\}, T, R_1 \cup R_2 \cup \{S \rightarrow S_1 S_2\})$ (where we assume that $S \notin N_1 \cup N_2$). (By the way, $L(G) = L(G_1)L(G_2)$.) If G_1 and G_2 are unambiguous, then so is G .
- (d) Suppose that for each $w \in T^*$ and $A \in N - \{S\}$ there exists at most one leftmost derivation which produces w from A . Suppose also that S does not appear on the right-hand side of any rule, namely for every $A \rightarrow \alpha \in R$ (with any $A \in N$) we have $\alpha \in (N \cup T - \{S\})^*$. Then G is unambiguous.

5. (a) Suppose G is defined by the rules

$$\begin{aligned} S &\rightarrow Aab \mid Aba \mid c, \\ A &\rightarrow SAc \mid AAb \mid ada, \end{aligned}$$

and we employ the algorithm discussed in class for eliminating left-recursion (with $A_1 = S, A_2 = A$). Then the grammar we obtain contains exactly eight rules.

- (b) Call a grammar *indirectly semi-recursive* (for the purposes of this question) if there exist $A \in N$ and $\alpha, \beta \in (N \cup T)^*$ such that $A \xrightarrow{+} \alpha A \beta$. If $L(G)$ is infinite then there does not exist a grammar, equivalent to G , that is not indirectly semi-recursive.
- (c) Denote by $l(G)$ the sum of lengths of the right-hand sides of all rules in R , namely:

$$l(G) = \sum_{A \rightarrow \alpha \in R} |\alpha|.$$

For a grammar G with neither ε -productions nor unit productions, let $\text{Chomsky}(G)$ be the grammar in Chomsky Normal Form, equivalent to G , constructed according to the algorithm presented in class. Then there exists a function $f : \mathbf{N} \rightarrow \mathbf{N}$ (where \mathbf{N} denotes the set of positive integers) such that $l(\text{Chomsky}(G)) \leq f(l(G))$ for every grammar with neither ε -productions nor unit productions nor useless letters.

- (d) Suppose all rules in R are of one of the three forms $A \rightarrow b$, $A \rightarrow bC$, and $A \rightarrow Bc$. Then a slight modification of the CYK algorithm yields a parsing algorithm that works in time $O(n^2)$ for words on length n .
6. (a) The number of stages (not including stage 0) in the algorithm presented in class for computing the set **FIRST** is at most $|T|$.
- (b) For $A \in N$, denote by **DIRECT_FOLLOW**(A) the set of all letters $X \in N \cup T$ for which R includes a rule of the form $B \rightarrow \alpha AX \beta$ for some $B \in N$ and $\alpha, \beta \in (N \cup T)^*$. Then:

$$\text{FOLLOW}(A) = \bigcup_{X \in \text{DIRECT_FOLLOW}(A)} \text{FIRST}(X).$$

- (c) If R includes the four rules $A \rightarrow aB$, $A \rightarrow \varepsilon$, $B \rightarrow aA$, and $B \rightarrow \varepsilon$ (in addition to other rules), then the grammar is not $LL(1)$.
- (d) If R includes the rules $A \rightarrow SBS$ and $A \rightarrow SbS$ (in addition to other rules), and $L(G)$ is infinite, then G is not $LL(k)$ for any $k \geq 1$.

Solutions

4. (a) Ambiguity means that there are words that can be produced by more than one parse tree. However, there may still be words produced by a single tree. For example, the grammar defined by the rules

$$\begin{aligned} S &\rightarrow A \mid a \mid b, \\ A &\rightarrow a, \end{aligned}$$

is ambiguous since the word a can be produced in two essentially different ways, namely $S \Rightarrow a$ and $S \Rightarrow A \Rightarrow a$. However, there is clearly a unique parse tree for the word b .

As a more interesting example, one may consider the grammar with **if-then** and **if-then-else**, discussed in class. The grammar is ambiguous, yet words without any occurrence of **else** are produced by a unique parse tree, as are words with an equal number of occurrences of **if** and **else**.

- (b) The classical construction of a grammar that accepts the language $L(G_1) \cup L(G_2)$ (obtained by an addition of a new start symbol S and the rules $S \rightarrow S_1$ and $S \rightarrow S_2$) is easily seen to provide an unambiguous grammar in our case.
- (c) Suppose G_1 and G_2 are defined by the rules

$$S_1 \rightarrow a \mid ab,$$

and

$$S_2 \rightarrow a \mid ba,$$

respectively. Then the word aba is produced in G in two ways: $S \Rightarrow S_1 S_2 \Rightarrow a S_2 \Rightarrow aba$, and $S \Rightarrow S_1 S_2 \Rightarrow ab S_2 \Rightarrow aba$. Thus, while both G_1 and G_2 are clearly unambiguous, G is ambiguous.

- (d) The grammar defined by the rules

$$\begin{aligned} S &\rightarrow A \mid a, \\ A &\rightarrow a, \end{aligned}$$

satisfies the required condition, yet it is obviously ambiguous as the word a may be obtained in two ways.

Thus, (a) and (b) are true.

5. (a) The rules for S do not have left-recursion. For A , we first need to replace the rule $A \rightarrow SAc$ by rules whose right-hand side does not start with the letter S . We obtain the following rules for A :

$$A \rightarrow AabAc \mid AbaAc \mid cAc \mid AAb \mid ada.$$

Now we need to get rid of the direct left-recursion we still have for A . We add a new non-terminal A' , and instead of the five above rules for A obtain the rules

$$\begin{aligned} A &\rightarrow cAcA' \mid adaA', \\ A' &\rightarrow abAcA' \mid baAcA' \mid AbA' \mid \varepsilon. \end{aligned}$$

The new grammar has altogether nine rules.

- (b) Consider any parse tree corresponding to a grammar that is not indirectly semi-recursive. Take any path from the root to a leaf of the tree. The assumed property of the grammar implies that all nodes along the way have distinct labels. Hence the height of the tree is bounded by $|N|$. It follows that $L(G)$ is finite.
- (c) In the process of passing from G to $\text{Chomsky}(G)$, we first add a rule $C_a \rightarrow a$ for each terminal a . Next we replace each rule of the form $A \rightarrow B_1B_2 \dots B_k$ with right-hand side of length $k \geq 3$ by $k - 1$ rules with right-hand sides of length 2 each. It follows that $l(\text{Chomsky}(G)) \leq |T| + 2l(G)$. Since all letters are useful, each terminal appears on the right-hand side of at least one rule, and consequently $|T|$ is bounded above by $l(G)$. It follows that $l(\text{Chomsky}(G)) \leq 3l(G)$, so that we may take f as the function given by $f(m) = 3m$.
- (d) We proceed as in the CYK algorithm. This time, the question for which non-terminals A and subwords $a_i a_{i+1} \dots a_j$ of the given input word $a_1 a_2 \dots a_n$ we have $A \xrightarrow{*} a_i a_{i+1} \dots a_j$ reduces to a bounded number of questions of the same form, with $a_i a_{i+1} \dots a_j$ replaced by one of the subwords $a_i a_{i+1} \dots a_{j-1}$ and $a_{i+1} a_{i+2} \dots a_j$, each of length $j - i$. Thus, for each $k \leq n$, we answer the above question for words of length k in time $O(k)$. Going over all k up to n , we complete the work in time $O(n^2)$.

Thus, (b), (c) and (d) are true.

6. (a) At the first stage, we find that a terminal a belongs to the **FIRST** set of a non-terminal A if $A \rightarrow \alpha a \beta \in R$ for some α, β with **Nullable**(α). At the second stage we find the same if $A \rightarrow \alpha B \beta \in R$, where **FIRST**(B) is known to include the terminal a from the preceding stage and **Nullable**(α), and so forth. It follows that the number of stages is bounded above by $|N|$. This bound cannot be reduced in general, as the grammar defined by the rules

$$\begin{aligned} S &\rightarrow A_1, \\ A_1 &\rightarrow A_2, \\ &\dots \\ A_{m-2} &\rightarrow A_{m-1}, \\ A_{m-1} &\rightarrow a \mid \varepsilon, \end{aligned}$$

shows.

- (b) In the proposed equality, indeed every letter belonging to the right-hand side belongs to the left-hand side too. In fact, let $a \in \mathbf{FIRST}(X)$, say $X \xRightarrow{*} \gamma a \delta$ for some $\gamma, \delta \in (N \cup T)^*$ with **Nullable**(γ), where $X \in \mathbf{DIRECT.FOLLOW}(A)$. Since all letters are useful, for suitable $\alpha', \beta' \in (N \cup T)^*$ we have

$$S \xRightarrow{*} \alpha' B \beta' \xRightarrow{*} \alpha' \alpha A X \beta \beta' \xRightarrow{*} \alpha' \alpha A \gamma a \delta \beta \beta' \xRightarrow{*} \alpha' \alpha A a \delta \beta \beta',$$

so that $a \in \mathbf{FOLLOW}(A)$.

However, the inverse inclusion is false in general. If **Nullable**(X), then elements of **FIRST**(β) also belong to **FOLLOW**(A). For example, consider the grammar defined by the rules

$$\begin{aligned} S &\rightarrow ABC, \\ A &\rightarrow a, \\ B &\rightarrow b \mid \varepsilon, \\ C &\rightarrow c. \end{aligned}$$

One verifies easily that $\mathbf{FOLLOW}(A) = \{b, c\}$, while

$$\bigcup_{X \in \mathbf{DIRECT.FOLLOW}(A)} \mathbf{FIRST}(X) = \{b\}.$$

- (c) It is easy to see that the grammar defined by the rules

$$\begin{aligned}
S &\rightarrow A, \\
A &\rightarrow aB \mid \varepsilon, \\
B &\rightarrow aA \mid \varepsilon.
\end{aligned}$$

is $LL(1)$.

- (d) Let $k \geq 1$ be arbitrary and fixed. We claim that G is not $LL(k)$. Since $L(G)$ is infinite, we may find a word $u \in L(G)$ with $|u| \geq k$. Consider a sequence of derivations of the form

$$S \xRightarrow{*} vA\beta \Longrightarrow vSBS\beta \xRightarrow{*} vuBS\beta \xRightarrow{*} vvw$$

for some $v, u, w \in T^*$ and $\beta \in (N \cup T)^*$. Now suppose we have to parse the input word vuw . After getting the sentential form $vA\beta$, the next k letters of the input that we have to match are the first k letters of u . Clearly, both rules $A \rightarrow SBS$ and $A \rightarrow SbS$ are equally suitable to be used at this stage (as far as the next k letters of the input are concerned). Hence G is not $LL(k)$.

Thus, only (d) is true.