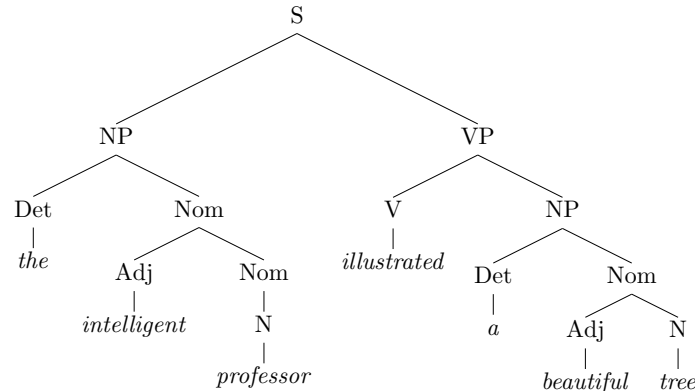


### 1.3 Ambiguity

If we are given a derivation of a word for some context-free grammar, we need not always represent it in a linear fashion like we did in the previous examples. We could alternatively represent it as a tree structure, where the root of the tree corresponds to the start nonterminal  $S$  and each branch of the tree adds a new nonterminal or terminal symbol. We refer to such trees as *parse trees*.

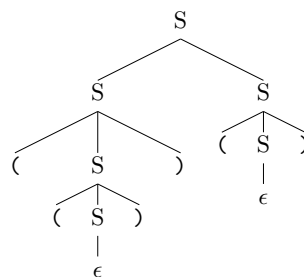
Parse trees are nothing new to linguists; the idea is used all the time to break down sentences or phrases into their constituent components, like nouns, verbs, and so on. In doing so, linguists are able to study the structures of sentences in different languages. For example, consider the following parse tree for an English sentence:



Here, the sentence ( $S$ ) is broken down into a noun phrase ( $NP$ ) and a verb phrase ( $VP$ ); the noun phrase is broken down further into a determiner ( $Det$ ) and a nominal ( $Nom$ ); and so on. There are all kinds of rules specifying exactly how we can break down English sentences in this way.

In every parse tree, the root of the tree is the start nonterminal  $S$ , the leaves of the tree contain terminal symbols from  $\Sigma$  (or  $\epsilon$ ), and all other vertices of the tree contain nonterminal symbols from  $V$ . If a parse tree contains an internal (non-leaf) vertex  $A$ , and all the children of the vertex  $A$  are labelled  $a_1, a_2, \dots, a_n$ , then the underlying grammar's rule set must contain a rule of the form  $A \rightarrow a_1 a_2 \dots a_n$ .

**Example 4.** Recall our derivation of the word  $(( )) ( )$  from the language of words with balanced parentheses in Example 3. We can represent the derivation of that word by the following parse tree:



For most grammars we deal with, there exists exactly one way to generate any given word in the language of the grammar. However, this is not always the case. There are some grammars that allow us to generate the same word in more than one way.

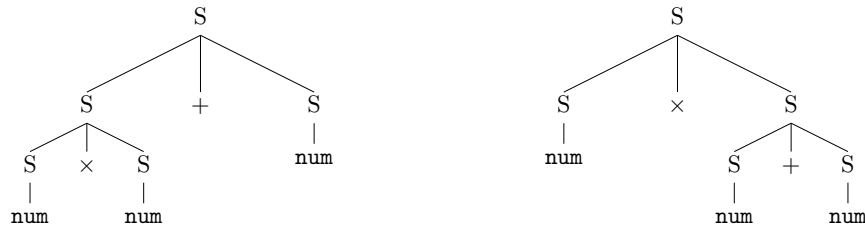
Perhaps one of the most well-known examples is the grammar generating the language of arithmetic expressions. If you recall grade school mathematics, you'll remember that there is an order of operations that specify the order in which we should apply arithmetic operations in a given expression.<sup>2</sup> We first evaluate expressions in parentheses, then exponents, then multiplications and divisions, and finally additions and subtractions.

<sup>2</sup>Sometimes referred to using the mnemonics BEDMAS, BIDMAS, BODMAS, or PEMDAS, depending on where you went to school.

Let's consider a simplified set of operations, where we only use parentheses, addition, and multiplication. The grammar generating the language of arithmetic expressions using these three operators together with the standard set of numbers is as follows:

$$\begin{aligned}
 S &\rightarrow (S) \\
 S &\rightarrow S + S \\
 S &\rightarrow S \times S \\
 S &\rightarrow \text{num}
 \end{aligned}$$

If we consider the expression  $\text{num} \times \text{num} + \text{num}$ , we discover that there exists more than one way to generate this expression, depending on whether we apply the  $+$  rule or the  $\times$  rule first. This is evidenced by the fact that there exist two parse trees for the same expression:

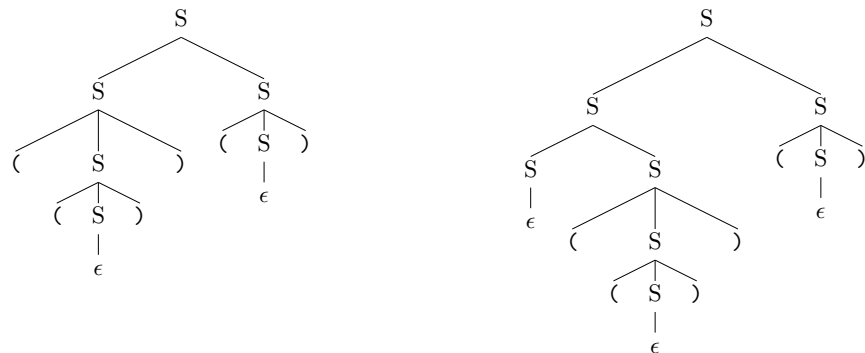


We don't need to do anything tricky in order to obtain these parse trees. In fact, both parse trees can be obtained simply by applying rules to each nonterminal from left to right; that is, at some level of the parse tree where there exists two nonterminals, we apply a rule to the first (left) nonterminal before the second (right) nonterminal. This process is known as a *leftmost derivation*.

If there exists some word in the language of a grammar for which there is more than one leftmost derivation of that word, then we say that the word is derived *ambiguously*. Likewise, the grammar producing that word is itself ambiguous.

**Definition 5** (Ambiguous context-free grammar). A context-free grammar  $G$  is ambiguous if there exists some word  $w \in L(G)$  that can be derived ambiguously.

**Example 6.** The grammar from Example 3 generating our language of words with balanced parentheses is ambiguous. Consider again the word  $(( )) ( )$ . There exist two different parse trees corresponding to leftmost derivations of this word; the left parse tree was given in Example 4 and the right parse tree is new:



## Removing Ambiguity

In some cases, if we have an ambiguous context-free grammar, we can create an equivalent context-free grammar with no ambiguity. For example, recall our three-operation arithmetic grammar from earlier:

$$\begin{aligned} S &\rightarrow (S) \\ S &\rightarrow S + S \\ S &\rightarrow S \times S \\ S &\rightarrow \text{num} \end{aligned}$$

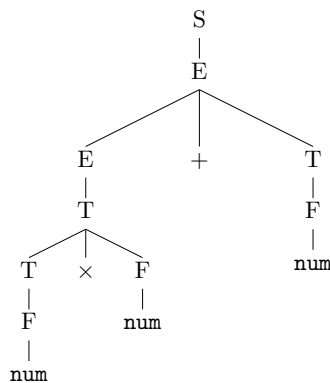
Nothing in this grammar forces us to use one rule before another, so we end up being able to derive the same word via different sequences of rule applications. However, we can construct an unambiguous grammar simply by adding a little more structure to our rules—specifically, by adding a few more nonterminals:

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T \times F \mid F \\ F &\rightarrow (E) \mid \text{num} \end{aligned}$$

Now, each nonterminal plays a particular role. The nonterminal  $S$ , as usual, serves as our starting point and produces an expression,  $E$ . Each expression consists of subexpressions  $E$  or additive terms  $T$ . Likewise, each term consists of subterms  $T$  or multiplicative factors  $F$ . Finally, each factor can be either a `num` or a parenthesized subexpression, starting the whole process over again.

We won't verify here that this grammar is in fact equivalent to our original one, though we can intuit that they each generate the same language. Suffice it to say that, with this revised grammar, we're able to ensure that the addition rule is always applied before the multiplication rule.

Our revised grammar now allows us to draw one unambiguous parse tree for the expression `num × num + num`:



## Inherent Ambiguity

Unfortunately, there is no general procedure or algorithm for removing ambiguity from a context-free grammar; indeed, it isn't even possible to remove ambiguity in some cases. Some context-free languages are *inherently ambiguous*, meaning that any grammar generating the language will have some unavoidable ambiguous component to it.

**Example 7.** Consider the language  $L_{\text{twoequal}} = \{a^i b^j c^k \mid i, j, k \geq 0 \text{ and } i = j \text{ or } j = k\}$  over the alphabet  $\Sigma = \{a, b, c\}$ . This language contains all words that have either the same number of `a`s and `b`s or the same number of `b`s and `c`s.

We can generate this language using the following grammar:

$$\begin{aligned} S &\rightarrow S_1 \mid S_2 \\ S_1 &\rightarrow S_1 c \mid A \\ A &\rightarrow aAb \mid \epsilon \\ S_2 &\rightarrow aS_2 \mid B \\ B &\rightarrow bBc \mid \epsilon \end{aligned}$$

The rules  $S_1$  and  $A$  generate words of the form  $a^n b^n c^m$  and the rules  $S_2$  and  $B$  generate words of the form  $a^m b^n c^n$ , each where  $m, n \geq 0$ .

Now, consider words of the form  $a^n b^n c^n$ , where  $n \geq 0$ . All words of this form belong to the language  $L_{\text{twoequal}}$ , but each such word has two distinct derivations in this grammar: it can be generated either by the rules  $S_1$  and  $A$ , or by the rules  $S_2$  and  $B$ .

While the formal proof that this language is inherently ambiguous is quite long, we can intuitively see that (for instance) any grammar generating this language must have rules similar to  $S_1$  and  $A$  to produce balanced pairs of  $a$ s and  $b$ s followed by some number of  $c$ s. We can make a similar argument for the rules  $S_2$  and  $B$ . Thus, any grammar for this language will include some degree of ambiguity.

## 1.4 Normal Forms

Up to now, we've imposed no restrictions on the form of each rule in our context-free grammars. As long as each rule of our context-free grammar looked like  $A \rightarrow \alpha$ , where  $A$  is a nonterminal symbol and  $\alpha$  is a combination of terminal and nonterminal symbols, we were happy.

However, computers (and, by extension, the people who program computers) like having structure. For instance, a compiler for a programming language usually incorporates a context-free grammar into its workflow at some point during the compilation of a program, and having a highly-structured grammar makes the compiler's job of determining which rule to apply both easier and faster.

Therefore, at times, we might like to transform a context-free grammar into a *normal form*; that is, to modify the grammar in such a way that each rule of the grammar follows a canonical form or template. There are many normal forms to choose from, and each one comes with its own benefits. However, here we will focus our attention on arguably the most popular normal form.

### Chomsky Normal Form

The *Chomsky normal form*, as the name suggests, was first studied by Noam Chomsky in the late 1950s as he attempted to develop a model of natural language using grammars. A grammar in Chomsky normal form is one where each rule either has two nonterminal symbols or one terminal symbol on the right-hand side.

**Definition 8** (Chomsky normal form). A context-free grammar is in Chomsky normal form if every rule in the grammar is of one of the two following forms:

1.  $A \rightarrow BC$  for  $A, B, C \in V$  with  $B, C \neq S$ ; or
2.  $A \rightarrow a$  for  $A \in V$  and  $a \in \Sigma$ .

Additionally, we may allow the rule  $S \rightarrow \epsilon$ .

The main benefit of converting a grammar into Chomsky normal form comes in how we can represent and store derivations of words in memory. Since each rule derives either two nonterminal symbols or one terminal symbol, every parse tree will have a branching factor of either 2 or 1. This fact allows us to use efficient data structures for representing binary trees in memory, as well as to apply efficient algorithms to process parse trees and derivations. Moreover, the number of steps in a derivation using a grammar in Chomsky normal form is easy to bound: if the grammar generates a word  $w$ , then the derivation of  $w$  will contain  $|w| - 1$  applications of a rule of the first form and  $|w|$  applications of a rule of the second form.

**Example 9.** Let  $\Sigma = \{a, b\}$ , and consider the following two grammars. Each grammar generates words consisting of one  $b$  surrounded on either side by zero or more  $a$ s. The grammar on the left is not in Chomsky normal form. The grammar on the right is in Chomsky normal form, and it is equivalent to the grammar on the left.

$$\begin{array}{ll} S \rightarrow AbA & S_0 \rightarrow TA \mid BA \mid AB \mid b \\ A \rightarrow Aa \mid \epsilon & T \rightarrow AB \\ & A \rightarrow AC \mid a \\ & B \rightarrow b \\ & C \rightarrow a \end{array}$$

Every context-free grammar can be converted into a context-free grammar in Chomsky normal form, and the conversion process consists of four steps:

1. **START:** Replace the start nonterminal.

Add a new start nonterminal  $S_0$  together with a new rule  $S_0 \rightarrow S$ , where  $S$  is the start nonterminal of the original grammar. This guarantees that  $S_0$  will not occur on the right-hand side of any rule.

2. **EPS:** Remove epsilon rules.

Remove all rules of the form  $A \rightarrow \epsilon$ , where  $A \neq S$ . For each occurrence of  $A$  on the right-hand side of a rule in the original grammar, add a new rule with that occurrence of  $A$  removed from the right-hand side; that is, convert all rules of the form  $X \rightarrow \alpha A \beta$  to  $X \rightarrow \alpha \beta$ , where  $\alpha$  and  $\beta$  are combinations of nonterminal and terminal symbols. Note that this applies to each *occurrence* of  $A$ , so a rule of the form  $X \rightarrow \alpha A \beta A \gamma$  would contribute three new rules:  $X \rightarrow \alpha \beta A \gamma$ ,  $X \rightarrow \alpha A \beta \gamma$ , and  $X \rightarrow \alpha \beta \gamma$ .

If there exists a rule of the form  $X \rightarrow A$ , then add a new rule of the form  $X \rightarrow \epsilon$  unless we previously removed a rule of that form.

Repeat until all epsilon rules not involving the original start nonterminal are removed.

3. **UNIT:** Remove unit rules.

Unit rules are rules of the form  $A \rightarrow B$ , where  $A$  and  $B$  are nonterminal symbols. Remove all rules of this form.

If there exists a rule of the form  $B \rightarrow \alpha$ , where  $\alpha$  is a combination of nonterminal and terminal symbols, then add a new rule  $A \rightarrow \alpha$  unless we previously removed a unit rule of that form.

Repeat until all unit rules are removed.

4. **BIN:** Remove rules with more than two nonterminal or terminal symbols on the right-hand side.

Replace each rule of the form  $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_k$ , where  $k \geq 3$  and each  $\alpha_i$  is either a nonterminal or terminal symbol, with a series of new rules  $A \rightarrow \alpha_1 A_1$ ,  $A_1 \rightarrow \alpha_2 A_2$ ,  $\dots$ ,  $A_{k-2} \rightarrow \alpha_{k-1} \alpha_k$ . Each  $A_i$  is a new nonterminal symbol.

If  $\alpha_i$  is a terminal symbol in one of our new rules  $A_{i-1} \rightarrow \alpha_i A_i$ , then remove this rule, replace  $\alpha_i$  with a new nonterminal symbol  $B_i$ , and add two new rules  $B_i \rightarrow \alpha_i$  and  $A_{i-1} \rightarrow B_i A_i$ .

**Example 10.** Consider the following grammar not in Chomsky normal form:

$$\begin{array}{l} S \rightarrow ASB \\ A \rightarrow aAS \mid a \mid \epsilon \\ B \rightarrow SbS \mid A \mid bb \end{array}$$

We will convert this grammar to an equivalent grammar in Chomsky normal form.

1. **START:** Replace the start nonterminal.

Adding a new start nonterminal and the rule  $S_0 \rightarrow S$  gives us the following:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \\ A &\rightarrow aAS \mid a \mid \epsilon \\ B &\rightarrow SbS \mid A \mid bb \end{aligned}$$

2. **EPS:** Remove epsilon rules.

The first epsilon rule we will remove is  $A \rightarrow \epsilon$ . Since  $A$  occurs on the right-hand side of two other rules ( $S \rightarrow ASB$  and  $A \rightarrow aAS$ ), removing this rule means we must add two new rules  $S \rightarrow SB$  and  $A \rightarrow aS$ . Additionally, since we have a rule  $B \rightarrow A$ , we must add a new rule  $B \rightarrow \epsilon$ . This gives the following:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \mid SB \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid A \mid bb \mid \epsilon \end{aligned}$$

Next, we remove the epsilon rule  $B \rightarrow \epsilon$  that we just added. Since  $B$  occurs on the right-hand side of two other rules ( $S \rightarrow ASB$  and  $S \rightarrow SB$ ), removing this rule means we must add two new rules  $S \rightarrow AS$  and  $S \rightarrow S$ . This gives the following:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \mid SB \mid AS \mid S \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid A \mid bb \end{aligned}$$

3. **UNIT:** Remove unit rules.

We first remove the unit rule  $B \rightarrow A$ . We do so by replacing  $A$  on the right-hand side of the rule with everything that can be produced by a rule of the form  $A \rightarrow \alpha$ :

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \mid SB \mid AS \mid S \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

Next, we remove the unit rule  $S \rightarrow S$ . This change is more straightforward, as we don't need to modify the  $S$  rule in any other way:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \mid SB \mid AS \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

Lastly, we remove the unit rule  $S_0 \rightarrow S$ . Again, we replace  $S$  on the right-hand side of the rule with everything that can be produced by a rule of the form  $S \rightarrow \alpha$ :

$$\begin{aligned} S_0 &\rightarrow ASB \mid SB \mid AS \\ S &\rightarrow ASB \mid SB \mid AS \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

4. **BIN**: Remove rules with more than two nonterminal or terminal symbols on the right-hand side.

Starting from the top, we replace the rule  $S_0 \rightarrow ASB$  with the rules  $S_0 \rightarrow AU_1$  and  $U_1 \rightarrow SB$ :

$$\begin{aligned} S_0 &\rightarrow AU_1 \mid SB \mid AS & U_1 &\rightarrow SB \\ S &\rightarrow ASB \mid SB \mid AS \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

Similarly, we replace the rule  $S \rightarrow ASB$  with the rules  $S \rightarrow AU_2$  and  $U_2 \rightarrow SB$ :

$$\begin{aligned} S_0 &\rightarrow AU_1 \mid SB \mid AS & U_1 &\rightarrow SB \\ S &\rightarrow AU_2 \mid SB \mid AS & U_2 &\rightarrow SB \\ A &\rightarrow aAS \mid a \mid aS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

Next, we replace the rule  $A \rightarrow aAS$  with the rules  $A \rightarrow aU_3$  and  $U_3 \rightarrow AS$ :

$$\begin{aligned} S_0 &\rightarrow AU_1 \mid SB \mid AS & U_1 &\rightarrow SB \\ S &\rightarrow AU_2 \mid SB \mid AS & U_2 &\rightarrow SB \\ A &\rightarrow aU_3 \mid a \mid aS & U_3 &\rightarrow AS \\ B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \end{aligned}$$

Moving along, we replace the rules  $B \rightarrow SbS$  and  $B \rightarrow aAS$  with the set of rules  $B \rightarrow SU_4$ ,  $B \rightarrow aU_5$ ,  $U_4 \rightarrow bS$ , and  $U_5 \rightarrow AS$ :

$$\begin{aligned} S_0 &\rightarrow AU_1 \mid SB \mid AS & U_1 &\rightarrow SB \\ S &\rightarrow AU_2 \mid SB \mid AS & U_2 &\rightarrow SB \\ A &\rightarrow aU_3 \mid a \mid aS & U_3 &\rightarrow AS \\ B &\rightarrow SU_4 \mid bb \mid aU_5 \mid a \mid aS & U_4 &\rightarrow bS \\ & & U_5 &\rightarrow AS \end{aligned}$$

Lastly, we replace all rules containing one terminal and one nonterminal symbol on the right-hand side. We introduce two new rules  $V_1 \rightarrow a$  and  $V_2 \rightarrow b$  and make the appropriate changes to other rules:

$$\begin{aligned} S_0 &\rightarrow AU_1 \mid SB \mid AS & U_1 &\rightarrow SB & V_1 &\rightarrow a \\ S &\rightarrow AU_2 \mid SB \mid AS & U_2 &\rightarrow SB & V_2 &\rightarrow b \\ A &\rightarrow V_1U_3 \mid a \mid V_1S & U_3 &\rightarrow AS \\ B &\rightarrow SU_4 \mid V_2V_2 \mid V_1U_5 \mid a \mid V_1S & U_4 &\rightarrow V_2S \\ & & U_5 &\rightarrow AS \end{aligned}$$