# 1   Grammars and Context-Free Languages

Recall that, in our discussion on regular languages, we introduced the notion of a regular expression. This expression essentially performed a kind of pattern matching to accept words of a certain form and reject words not of that form.

We can take this rough idea of matching patterns in words and modify it to work not just for symbols within each word, but for the structure and composition of the word itself. We are able to do so using the notion of *grammars*, which essentially provide a set of *rules* that we can follow to produce words that belong to a certain language. (Note that, while we used regular expressions to match a *given* word, we use grammars to produce a *new* word.) If a grammar produces words that belong to a certain language, then we say the grammar *generates* that language.

The idea of using grammars with languages is nothing new; linguists have been using grammars to study natural languages for centuries, dating as far back as the fourth century BCE with the work of the Indian grammarian Pāṇini. Only with the advent of computer science itself has the notion of grammars been applied to formal languages and programming languages, starting with the work of American linguist Noam Chomsky in the 1950s.

While we didn't mention it in the previous lecture, there exists the notion of a *regular grammar*, which corresponds exactly to the class of regular languages. The rules of a regular grammar take on one of the following forms:[1]

1. $A \to \epsilon$;

2. $A \to a$ for some $a \in \Sigma$; and

3. $A \to aB$ for some $a \in \Sigma$.

In each rule, the lowercase letters correspond to alphabet symbols and the uppercase letters correspond to other rules that we can use to produce more symbols of the word. In particular, though, observe how we can associate each rule with the action a finite automaton takes when reading an input word: Rule 1 corresponds to reading no symbol (e.g., while following an epsilon transition), Rule 2 corresponds to reading the last symbol of the word, and Rule 3 corresponds to reading some symbol of the word and transitioning to another state to read the next symbol.

However, this lecture isn't about regular grammars! We already have a number of ways to represent regular languages, and we know also that there exist some non-regular languages. Thus, instead of dwelling on the regular languages, here we will take our first look at a larger class: the class of *context-free languages*.

## 1.1   Context-Free Grammars

If you look at the specification manual for any programming language, you will likely find tucked away somewhere in the documentation a grammar for that language. This grammar, which could number into the tens of pages, describes precisely what the structure of a program written in that language should look like. Indeed, this grammar is exactly what the language compiler uses to check for syntax errors in your program!

---

[1]Strictly speaking, these rules define the class of *right-linear* regular grammars. A similar set of rules exists for *left-linear* regular grammars, but both types of grammars generate the regular languages.

As an example, let's consider one small excerpt from the grammar given in the third edition of the *Java Language Specification*:

```
Statement:
    Block
    assert Expression [ : Expression] ;
    if ParExpression Statement [else Statement]
    for ( ForControl ) Statement
    while ParExpression Statement
    do Statement while ParExpression ;
    try Block ( Catches | [Catches] finally Block )
    switch ParExpression { SwitchBlockStatementGroups }
    synchronized ParExpression Block
    return [Expression] ;
    throw Expression ;
    break [Identifier]
    continue [Identifier]
    ;
    StatementExpression ;
    Identifier   :   Statement
```

This part of the Java grammar checks statement blocks such as assignments, if-else blocks, for loops, and so on. All of the words written with an `Uppercase` letter or written in `CamelCase` correspond to rules, and all of the words written in `lowercase` correspond to language keywords. For example, the `if` rule on the fourth line checks that every if-else block in a program conforms to the syntax that the compiler expects: the block begins with the keyword `if` together with some parenthesized expression, followed by some statement or sequence of instructions, and ending with an optional `else` block.

This Java grammar is an example of a *context-free grammar*. Like the regular grammars we saw earlier, a context-free grammar consists of a set of rules that we can use to generate valid programs in Java. Unlike regular grammars, however, these rules take on a much more general form: for example, we can replace the `Statement` rule with any combination of keywords and other rules, as specified by that part of the grammar.

Before we look at some other examples, let's formalize the notion of a context-free grammar.

**Definition 1** (Context-free grammar)**.** A context-free grammar is a tuple $(V, \Sigma, R, S)$, where

- $V$ is a finite set of elements called *nonterminal symbols*;

- $\Sigma$ is a finite set of elements called *terminal symbols*, where $\Sigma \cap V = \emptyset$;

- $R$ is a finite set of *rules*, where each rule consists of a nonterminal on the left-hand side and a combination of nonterminals and terminals on the right-hand side; and

- $S \in V$ is the *start nonterminal*.

In a context-free grammar, the set of nonterminal symbols $V$ correspond to parts of a word that we have yet to "fill in" with terminal symbols from $\Sigma$. The set of rules $R$ tell us how we can perform this "filling in". If we have a rule of the form $A \to \alpha$, then we can replace any instance of $A$ in our word with the symbol $\alpha$. The start nonterminal $S$ is self-explanatory; it is the first thing in our word that we "fill in".

Returning to our Java grammar example, we can see that (for example) some of the nonterminals in the grammar include `Statement`, `Block`, `Identifier`, and `ParExpression`, while some of the terminals include `if`, `while`, `for`, and `;` (semicolon).

Importantly, we have in our definition of a context-free grammar that $\Sigma \cap V = \emptyset$; that is, the set of terminals and the set of nonterminals must be disjoint. This is to prevent the grammar from confusing terminals and nonterminals. (Incidentally, this is why the Java language designers used uppercase letters in their nonterminals and lowercase letters in their terminals.)

The sequence of rule applications we follow beginning with the start nonterminal $S$ and ending with a completed word containing symbols from $\Sigma$ is called a *derivation*. Each word of the form $(V \cup \Sigma)^*$ produced during a derivation is sometimes referred to as a *sentential form*.

For any nonterminal $A$ and terminals $u$, $w$, and $v$, if we have a rule $A \to w$ in our grammar and some step of our derivation takes us from $uAv$ to $uwv$, then we say that $uAv$ *yields* $uwv$ and we write $uAv \Rightarrow uwv$. We can represent a sequence of "yields" relations using similar notation; given words $x$ and $y$, if $x = y$ or if there exists a sequence $x_1, x_2, \ldots, x_k$ where $k \geq 0$ such that

$$x \Rightarrow x_1 \Rightarrow x_2 \Rightarrow \cdots \Rightarrow x_k \Rightarrow y,$$

then we write $x \Rightarrow^* y$. (This is very similar to the Kleene star notation, where the star indicates zero or more "yields" relations taking us from $x$ to $y$.)

**Example 2.** Consider the context-free grammar where $V = \{S, A\}$, $\Sigma = \{\mathsf{a}, \mathsf{b}\}$, and $R$ contains two rules:

$$S \to \mathsf{a}A\mathsf{b}$$
$$A \to \mathsf{a}A\mathsf{b} \mid \epsilon$$

Using this context-free grammar, we can generate words like

$$S \Rightarrow \mathsf{a}A\mathsf{b} \Rightarrow \mathsf{a}\,\epsilon\,\mathsf{b} = \mathsf{ab},$$

$$S \Rightarrow \mathsf{a}A\mathsf{b} \Rightarrow \mathsf{a}\,\mathsf{a}A\mathsf{b}\,\mathsf{b} \Rightarrow \mathsf{aa}\,\epsilon\,\mathsf{bb} = \mathsf{aabb}, \text{ and}$$

$$S \Rightarrow \mathsf{a}A\mathsf{b} \Rightarrow \mathsf{a}\,\mathsf{a}A\mathsf{b}\,\mathsf{b} \Rightarrow \mathsf{aa}\,\mathsf{a}A\mathsf{b}\,\mathsf{bb} \Rightarrow \mathsf{aaa}\,\epsilon\,\mathsf{bbb} = \mathsf{aaabbb},$$

and so on. For each step, the highlighted symbols indicate which symbols were added at that step. As we can see, this context-free grammar generates all words over the alphabet $\Sigma = \{\mathsf{a}, \mathsf{b}\}$ where the number of $\mathsf{a}$s is equal to the number of $\mathsf{b}$s, there is at least one $\mathsf{a}$ and one $\mathsf{b}$, and all $\mathsf{a}$s come before $\mathsf{b}$s in the word.

Observe that our rule $A$ in Example 2 included a vertical bar. This is simply a shorthand for writing multiple rules where each rule contains $A$ on the left-hand side. Writing $A \to \mathsf{a}A\mathsf{b} \mid \epsilon$ is therefore equivalent to writing

$$A \to \mathsf{a}A\mathsf{b}$$
$$A \to \epsilon$$

There are very few limitations we must abide by when we write rules for a context-free grammar. All we need to ensure is that the left-hand side of each rule consists of exactly one nonterminal by itself. The right-hand side of each rule can contain any combination of terminals and nonterminals, including the empty word $\epsilon$.

**Example 3.** Consider the context-free grammar where $V = \{S\}$, $\Sigma = \{(,)\}$, and $R$ contains one rule:

$$S \to (S) \mid SS \mid \epsilon$$

This rule allows us to surround an occurrence of $S$ with parentheses, to "duplicate" an occurrence of $S$, or to replace some occurrence of $S$ with $\epsilon$, effectively removing that occurrence of $S$ from the derivation.

Using this context-free grammar, we can generate a word like

$$S \Rightarrow SS \Rightarrow (S)\,S \Rightarrow (\,(S)\,)S \Rightarrow ((\,\epsilon\,))S \Rightarrow (())S \Rightarrow (())\,(S) \Rightarrow (())(\,\epsilon\,) = (())().$$

Again, the highlighted symbols indicate which symbols were added at a given step. This context-free grammar generates all words over the alphabet $\Sigma = \{(,)\}$ where each word contains *balanced parentheses*: every opening parenthesis is matched by a closing parenthesis, and each pair of parentheses is correctly nested.
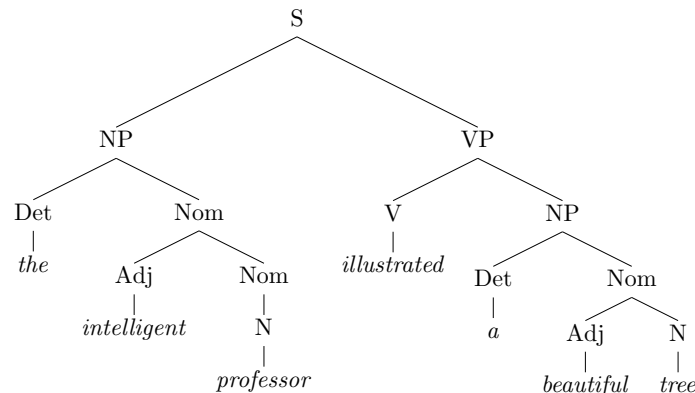
We can define the *language of a grammar G* over an alphabet $\Sigma$ by $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$. Thus, in Example 2, the language of the grammar is $L_{\mathtt{a=b}} = \{\mathtt{a}^n\mathtt{b}^n \mid n \geq 1\}$. Likewise, in Example 3, the language of the grammar is $L_{()} = \{w \in \{\mathtt{(},\mathtt{)}\}^* \mid$ all prefixes of $w$ contain no more $)$s than $($s, and $|w|_{(} = |w|_{)}\}$.[2]

Finally, in case you're wondering why we refer to these grammars as being *context-free*: the name stems from the fact that, given a rule of the form $A \to \alpha$, we can replace any occurrence of $A$ with $\alpha$ without looking at the *context* around $A$; that is, without looking at the symbols to the left and right of $A$. Thus, the grammar is free of context when we apply a given rule.

## 1.2 Ambiguous Context-Free Grammars

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 visually as a tree structure, where the root of the tree corresponds to the starting nonterminal $S$ and each branch of the tree adds a new nonterminal or terminal symbol. We refer to such trees as *parse trees*.
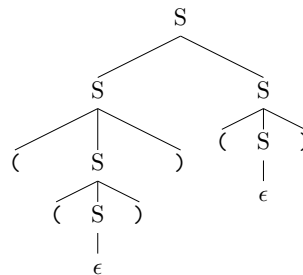
Parse trees are nothing new to linguists; the idea is used all the time to break down sentences or phrases into their constituent components, like nouns, verbs, and so on. In doing so, linguists are able to study the structures of sentences in different languages. For example, consider the following 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.

Of course, just like we can use grammars with formal languages, we can use parse trees to represent the derivation of any word in the language of a grammar. In our parse trees, the root of the tree is the starting 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 $a_1, a_2, \ldots, a_n$, then the underlying grammar's rule set contains a rule of the form $A \to a_1 a_2 \ldots a_n$.

**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:
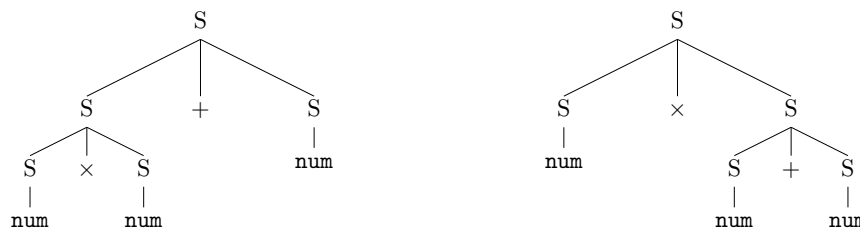


---

[2]The language of balanced parentheses is also known as the *Dyck language*.

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 wherein the same word can be generated 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.[3] We first evaluate expressions in parentheses, then exponents, then multiplications and divisions, and finally additions and subtractions.

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:

$$S \to (S)$$
$$S \to S + S$$
$$S \to S \times S$$
$$S \to \texttt{num}$$

If we consider the expression $\texttt{num} \times \texttt{num} + \texttt{num}$, 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:
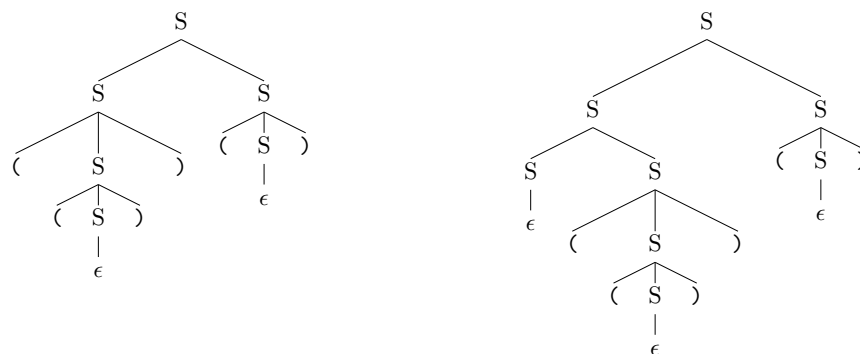


We also assume that, in both of these parse trees, we apply rules to each nonterminal from left to right; that is, at some level of the parse tree where there exists two nonterminals, we apply a rule to the first (left) nonterminal before the second (right) nonterminal. This process is known as a *leftmost 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*. This notion leads us to our main definition.

**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:
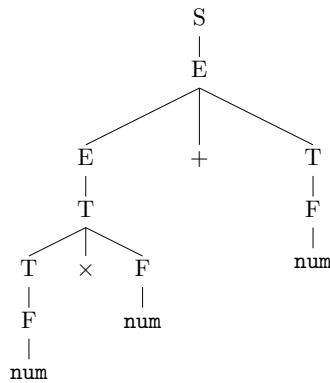


---

[3]Sometimes referred to using the mnemonics BEDMAS, BIDMAS, BODMAS, or PEMDAS, depending on where you went to school.

In some cases, if we have an ambiguous context-free grammar, we can create an equivalent context-free grammar with reduced or no ambiguity. For example, recalling our three-operation arithmetic grammar from earlier, we can construct an unambiguous grammar generating the same language of arithmetic expressions simply by adding a few "more structured" rules:

$$S \to E$$
$$E \to E + T \mid T$$
$$T \to T \times F \mid F$$
$$F \to (E) \mid \texttt{num}$$

With this grammar, we're able to ensure that the addition rule is applied before the multiplication rule. This allows us to draw a single, unambiguous parse tree for the expression $\texttt{num} \times \texttt{num} + \texttt{num}$:



However, this process of reducing or removing ambiguity cannot always be done. Some context-free grammars are *inherently ambiguous*, meaning that any grammar generating the specified language has some unavoidable ambiguous component to it, and this ambiguity cannot be removed.

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

We can generate this language using the following grammar:

$$S \to S_1 \mid S_2$$
$$S_1 \to S_1 \texttt{c} \mid A$$
$$A \to \texttt{a}A\texttt{b} \mid \epsilon$$
$$S_2 \to \texttt{a}S_2 \mid B$$
$$B \to \texttt{b}B\texttt{c} \mid \epsilon$$

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

Now, consider words of the form $\texttt{a}^n \texttt{b}^n \texttt{c}^n$, 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 grammar 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 $\texttt{a}$s and $\texttt{b}$s followed by some number of $\texttt{c}$s. We can make a similar argument for the rules $S_2$ and $B$. Thus, any grammar for this language has inherent ambiguity.

## 1.3    Normal Forms of Context-Free Grammars

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 $A \to \alpha$, 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, in some cases, 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 follows a canonical form or template.

The normal form we will focus on, the *Chomsky normal form*, was (as the name suggests) 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 such that 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. $A \to BC$ for $A, B, C \in V$ with $B, C \neq S$; or

2. $A \to a$ for $A \in V$ and $a \in \Sigma$.

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

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 $\Sigma = \{\mathtt{a}, \mathtt{b}\}$, and consider the following two grammars. Each grammar generates words consisting of one $\mathtt{b}$ surrounded on either side by zero or more $\mathtt{a}$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.

$$
\begin{aligned}
S &\to A\mathtt{b}A & S_0 &\to TA \mid BA \mid AB \mid \mathtt{b} \\
A &\to A\mathtt{a} \mid \epsilon & T &\to AB \\
& & A &\to AC \mid \mathtt{a} \\
& & B &\to \mathtt{b} \\
& & C &\to \mathtt{a}
\end{aligned}
$$

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 $S_0 \to S$, 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 $A \to \epsilon$, 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 $X \to \alpha A \beta$ to $X \to \alpha \beta$, 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 $X \to \alpha A \beta A \gamma$ would contribute three new rules: $X \to \alpha \beta A \gamma$, $X \to \alpha A \beta \gamma$, and $X \to \alpha \beta \gamma$.

If there exists a rule of the form $X \to A$, then add a new rule of the form $X \to \epsilon$ 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 $A \to B$, where $A$ and $B$ are nonterminal symbols. Remove all rules of this form.

If there exists a rule of the form $B \to \alpha$, where $\alpha$ is a combination of nonterminal and terminal symbols, then add a new rule $A \to \alpha$ 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 $A \to \alpha_1 \alpha_2 \ldots \alpha_k$, where $k \geq 3$ and each $\alpha_i$ is either a nonterminal or terminal symbol, with a series of new rules $A \to \alpha_1 A_1$, $A_1 \to \alpha_2 A_2$, $\ldots$, $A_{k-2} \to \alpha_{k-1} \alpha_k$. Each $A_i$ is a new nonterminal symbol.

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

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

$$S \to ASB$$
$$A \to \mathtt{a}AS \mid \mathtt{a} \mid \epsilon$$
$$B \to S\mathtt{b}S \mid A \mid \mathtt{bb}$$

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 $S_0 \to S$ gives us the following:

$$S_0 \to S$$
$$S \to ASB$$
$$A \to \mathtt{a}AS \mid \mathtt{a} \mid \epsilon$$
$$B \to S\mathtt{b}S \mid A \mid \mathtt{bb}$$

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

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

$$S_0 \to S$$
$$S \to ASB \mid SB$$
$$A \to \mathtt{a}AS \mid \mathtt{a} \mid \mathtt{a}S$$
$$B \to S\mathtt{b}S \mid A \mid \mathtt{bb} \mid \epsilon$$

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

$$S_0 \to S$$
$$S \to ASB \mid SB \mid AS \mid S$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid A \mid \mathsf{bb}$$

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

We first remove the unit rule $B \to A$. 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 $A \to \alpha$:

$$S_0 \to S$$
$$S \to ASB \mid SB \mid AS \mid S$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$

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

$$S_0 \to S$$
$$S \to ASB \mid SB \mid AS$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$

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

$$S_0 \to ASB \mid SB \mid AS$$
$$S \to ASB \mid SB \mid AS$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$

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 $S_0 \to ASB$ with the rules $S_0 \to AU_1$ and $U_1 \to SB$:

$$S_0 \to AU_1 \mid SB \mid AS$$
$$S \to ASB \mid SB \mid AS$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$U_1 \to SB$$

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

$$S_0 \to AU_1 \mid SB \mid AS$$
$$S \to AU_2 \mid SB \mid AS$$
$$A \to \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$U_1 \to SB$$
$$U_2 \to SB$$

Next, we replace the rule $A \to \mathsf{a}AS$ with the rules $A \to \mathsf{a}U_3$ and $U_3 \to AS$:

$$S_0 \to AU_1 \mid SB \mid AS$$
$$S \to AU_2 \mid SB \mid AS$$
$$A \to \mathsf{a}U_3 \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to S\mathsf{b}S \mid \mathsf{bb} \mid \mathsf{a}AS \mid \mathsf{a} \mid \mathsf{a}S$$
$$U_1 \to SB$$
$$U_2 \to SB$$
$$U_3 \to AS$$

Moving along, we replace the rules $B \to S\mathsf{b}S$ and $B \to \mathsf{a}AS$ with the set of rules $B \to SU_4$, $B \to \mathsf{a}U_5$, $U_4 \to \mathsf{b}S$, and $U_5 \to AS$:

$$S_0 \to AU_1 \mid SB \mid AS$$
$$S \to AU_2 \mid SB \mid AS$$
$$A \to \mathsf{a}U_3 \mid \mathsf{a} \mid \mathsf{a}S$$
$$B \to SU_4 \mid \mathsf{bb} \mid \mathsf{a}U_5 \mid \mathsf{a} \mid \mathsf{a}S$$
$$U_1 \to SB$$
$$U_2 \to SB$$
$$U_3 \to AS$$
$$U_4 \to \mathsf{b}S$$
$$U_5 \to AS$$

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

$$S_0 \to AU_1 \mid SB \mid AS$$
$$S \to AU_2 \mid SB \mid AS$$
$$A \to V_1U_3 \mid \mathsf{a} \mid V_1S$$
$$B \to SU_4 \mid V_2V_2 \mid V_1U_5 \mid \mathsf{a} \mid V_1S$$
$$U_1 \to SB$$
$$U_2 \to SB$$
$$U_3 \to AS$$
$$U_4 \to V_2S \qquad\qquad V_1 \to \mathsf{a}$$
$$U_5 \to AS \qquad\qquad V_2 \to \mathsf{b}$$

## 1.4 Context-Free Languages

With our notion of a context-free grammar, it's easy for us to define a context-free language. Just like we defined regular languages in terms of regular operations or regular expressions, we have the following definition for context-free languages in terms of context-free grammars.

**Definition 11** (Context-free language)**.** If some language $L$ is generated by a context-free grammar, then $L$ is context-free.

Thus, each of the languages in this lecture for which we constructed context-free grammars—the language of words $\mathsf{a}^n\mathsf{b}^n$, the language of balanced parentheses, the language of words $\mathsf{a}^i\mathsf{b}^j\mathsf{c}^k$ where $i = j$ or $j = k$— are context-free languages. Our definition also tells us that a language is context-free if we can construct a context-free grammar generating words in that language. However, be careful: any such context-free grammar must generate *all* and *only all* the words in the given language.

As a shorthand, we denote the class of languages generated by a context-free grammar by $\mathsf{CFG}$.

**Example 12.** Consider the language $L = \{\mathsf{a}^{2i}\mathsf{b}^i\mathsf{c}^{j+2} \mid i, j \geq 0\}$ over the alphabet $\Sigma = \{\mathsf{a}, \mathsf{b}, \mathsf{c}\}$. Here, we can see that the counts of $\mathsf{a}$s and $\mathsf{b}$s are related, while the number of $\mathsf{c}$s is independent of the number of $\mathsf{a}$s and $\mathsf{b}$s.

Let's construct a context-free grammar generating words in $L$. Evidently, we will need two kinds of rules: one rule will generate the $\mathsf{a}$s and $\mathsf{b}$s together, while the other rule will generate the $\mathsf{c}$s. We can use the start nonterminal to apply these rules in the correct order. Our context-free grammar will therefore look like the following:

$$S \rightarrow UV$$
$$U \rightarrow \mathsf{aa}U\mathsf{b} \mid \epsilon$$
$$V \rightarrow \mathsf{c}V \mid \mathsf{cc}$$

Let's now take a look at each rule in turn. The first rule, $S$, ensures that we apply the $U$ rule before the $V$ rule. This in turn ensures that all $\mathsf{a}$s and $\mathsf{b}$s occur before the $\mathsf{c}$s in the generated word. The second rule, $U$, either recursively produces two $\mathsf{a}$s and one $\mathsf{b}$ or produces the empty word. This ensures that we maintain the correct count of $2i$ $\mathsf{a}$s and $i$ $\mathsf{b}$s. Finally, the third rule, $V$, either recursively produces one $\mathsf{c}$ or produces the symbols $\mathsf{cc}$. This ensures that we have exactly $j + 2$ $\mathsf{c}$s in our generated word.

**Example 13.** Recall our context-free language $L_{\mathsf{a=b}} = \{\mathsf{a}^n\mathsf{b}^n \mid n \geq 1\}$. Here, let's consider a more general language: $L_{\mathrm{mixed\mathsf{a=b}}} = \{w \in \{\mathsf{a}, \mathsf{b}\}^* \mid |w|_\mathsf{a} = |w|_\mathsf{b}\}$. Observe that the main difference with this language is that the order of $\mathsf{a}$s and $\mathsf{b}$s no longer matters; we just need the same count of $\mathsf{a}$s and $\mathsf{b}$s. Can we construct a context-free grammar for $L_{\mathrm{mixed\mathsf{a=b}}}$?

Since order no longer matters, we just need our context-free grammar to generate a pair of $\mathsf{a}$s and $\mathsf{b}$s each time we add terminal symbols. For this, we can use essentially the same rule as we used in our context-free grammar for $L_{\mathsf{a=b}}$: $S \rightarrow \mathsf{a}S\mathsf{b} \mid \mathsf{b}S\mathsf{a}$. We also need a rule that allows us to mix the order of $\mathsf{a}$s and $\mathsf{b}$s; for instance, to place two $\mathsf{a}$s or two $\mathsf{b}$s next to each other instead of strictly nesting pairs. For this, we can use a rule similar to one we included in our context-free grammar for $L_{()}$: $S \rightarrow SS$.

Thus, our context-free grammar will look like the following:

$$S \rightarrow SS \mid \mathsf{a}S\mathsf{b} \mid \mathsf{b}S\mathsf{a} \mid \epsilon$$

# 2 Pushdown Automata

When we first introduced finite automata as a computational model for regular languages, we emphasized the facts that finite automata have no memory, no storage, and no ability to return to a previously-read symbol. Naturally, these restrictions limited the kinds of languages the model is able to recognize, and we showed that such restrictions resulted in the model recognizing exactly the class of regular languages.