

Caveat lector: This is the first edition of this lecture note. Please send bug reports and suggestions to jeffe@illinois.edu.

Imagine a piano keyboard, eh, 88 keys, only 88 and yet, and yet, hundreds of new melodies, new tunes, new harmonies are being composed upon hundreds of different keyboards every day in Dorset alone. Our language, tiger, our language: hundreds of thousands of available words, frillions of legitimate new ideas, so that I can say the following sentence and be utterly sure that nobody has ever said it before in the history of human communication: "Hold the newsreader's nose squarely, waiter, or friendly milk will countermand my trousers." Perfectly ordinary words, but never before put in that precise order. A unique child delivered of a unique mother.

— Stephen Fry, *A Bit of Fry and Laurie*, Series 1, Episode 3 (1989)

5 Context-Free Languages and Grammars

5.1 Definitions

Intuitively, a language is regular if it can be built from individual strings by concatenation, union, and repetition. In this note, we consider a wider class of **context-free** languages, which are languages that can be built from individual strings by concatenation, union, and *recursion*.

Formally, a language is context-free if and only if it has a certain type of recursive description known as a **context-free grammar**, which is a structure with the following components:

- A finite set Σ , whose elements are called **symbols** or **terminals**.
- A finite set Γ disjoint from Σ , whose elements are called **non-terminals**.
- A finite set R of **production rules** of the form $A \rightarrow w$, where $A \in \Gamma$ is a non-terminal and $w \in (\Sigma \cup \Gamma)^*$ is a string of symbols and variables.
- A **starting** non-terminal, typically denoted S .

For example, the following eight production rules describe a context free grammar with terminals $\Sigma = \{0, 1\}$ and non-terminals $\Gamma = \{S, A, B\}$:

$$\begin{array}{llll} S \rightarrow A & A \rightarrow 0A & B \rightarrow B1 & C \rightarrow \varepsilon \\ S \rightarrow B & A \rightarrow 0C & B \rightarrow C1 & C \rightarrow 0C1 \end{array}$$

Normally we write grammars more compactly by combining the right sides of all rules for each non-terminal into one list, with alternatives separated by vertical bars.¹ For example, the previous grammar can be written more compactly as follows:

$$\begin{array}{l} S \rightarrow A \mid B \\ A \rightarrow 0A \mid 0C \\ B \rightarrow B1 \mid C1 \\ C \rightarrow \varepsilon \mid 0C1 \end{array}$$

For the rest of this lecture, I will *almost* always use the following notational conventions.

¹Yes, this means we now have *three* symbols \cup , $+$, and $|$ with exactly the same meaning. Sigh.

- Monospaced digits (0, 1, 2, ...) and symbols (\diamond , \$, #, •, ...) are explicit terminals.
- Early lower-case Latin letters (a, b, c, ...) represent unknown or arbitrary terminals in Σ .
- Upper-case Latin letters (A, B, C, ...) and the letter S represent non-terminals in Γ .
- Late lower-case Latin letters (... , w, x, y, z) represent strings in $(\Sigma \cup \Gamma)^*$, whose characters could be either terminals or non-terminals.

We can *apply* a production rule to a string in $(\Sigma \cup \Gamma)^*$ by replacing any instance of the non-terminal on the left of the rule with the string on the right. More formally, for any strings $x, y, z \in (\Sigma \cup \Gamma)^*$ and any non-terminal $A \in \Gamma$, applying the production rule $A \rightarrow y$ to the string xAz yields the string xyz . We use the notation $xAz \rightsquigarrow xyz$ to describe this application. For example, we can apply the rule $C \rightarrow \theta C 1$ to the string $\theta \theta C 1 B A C \theta$ in two different ways:

$$\theta \theta \underline{C} 1 B A C \theta \rightsquigarrow \theta \theta \underline{\theta C 1} 1 B A C \theta \quad \theta \theta C 1 B A \underline{C} \theta \rightsquigarrow \theta \theta C 1 B A \underline{\theta C 1} \theta$$

More generally, for any strings $x, z \in (\Sigma \cup \Gamma)^*$, we say that z *derives from* x , written $x \rightsquigarrow^* z$, if we can transform x into z by applying a finite sequence of production rules, or more formally, if either

- $x = z$, or
- $x \rightsquigarrow y$ and $y \rightsquigarrow^* z$ for some string $y \in (\Sigma \cup \Gamma)^*$.

Straightforward definition-chasing implies that, for any strings $w, x, y, z \in (\Sigma \cup \Gamma)^*$, if $x \rightsquigarrow^* y$, then $wxz \rightsquigarrow^* wyz$.

The *language* $L(w)$ of any string $w \in (\Sigma \cup \Gamma)^*$ is the set of all strings in Σ^* that derive from w :

$$L(w) := \{x \in \Sigma^* \mid w \rightsquigarrow^* x\}.$$

The language *generated by* a context-free grammar G , denoted $L(G)$, is the language of its starting non-terminal. Finally, a language is *context-free* if it is generated by some context-free grammar.

Context-free grammars are sometimes used to model natural languages. In this context, the symbols are *words*, and the strings in the languages are *sentences*. For example, the following grammar describes a simple subset of English sentences. (Here I diverge from the usual notation conventions. Strings in $\langle \rangle$ are non-terminals, and regular strings are terminals.)

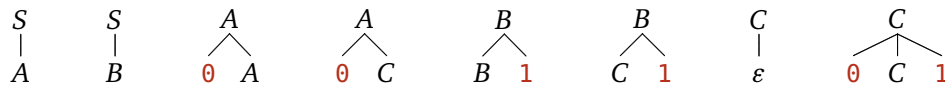
$$\begin{aligned} \langle \text{sentence} \rangle &\rightarrow \langle \text{noun phrase} \rangle \langle \text{verb phrase} \rangle \langle \text{noun phrase} \rangle \\ \langle \text{noun phrase} \rangle &\rightarrow \langle \text{adjective phrase} \rangle \langle \text{noun} \rangle \\ \langle \text{adj. phrase} \rangle &\rightarrow \langle \text{article} \rangle \mid \langle \text{possessive} \rangle \mid \langle \text{adjective phrase} \rangle \langle \text{adjective} \rangle \\ \langle \text{verb phrase} \rangle &\rightarrow \langle \text{verb} \rangle \mid \langle \text{adverb} \rangle \langle \text{verb phrase} \rangle \\ \langle \text{noun} \rangle &\rightarrow \text{dog} \mid \text{trousers} \mid \text{daughter} \mid \text{nose} \mid \text{homework} \mid \text{time lord} \mid \text{pony} \mid \dots \\ \langle \text{article} \rangle &\rightarrow \text{the} \mid \text{a} \mid \text{some} \mid \text{every} \mid \text{that} \mid \dots \\ \langle \text{possessive} \rangle &\rightarrow \langle \text{noun phrase} \rangle \text{'s} \mid \text{my} \mid \text{your} \mid \text{his} \mid \text{her} \mid \dots \\ \langle \text{adjective} \rangle &\rightarrow \text{friendly} \mid \text{furious} \mid \text{moist} \mid \text{green} \mid \text{severed} \mid \text{timey-wimey} \mid \text{little} \mid \dots \\ \langle \text{verb} \rangle &\rightarrow \text{ate} \mid \text{found} \mid \text{wrote} \mid \text{killed} \mid \text{mangled} \mid \text{saved} \mid \text{invented} \mid \text{broke} \mid \dots \\ \langle \text{adverb} \rangle &\rightarrow \text{squarely} \mid \text{incompetently} \mid \text{barely} \mid \text{sort of} \mid \text{awkwardly} \mid \text{totally} \mid \dots \end{aligned}$$

5.2 Parse Trees

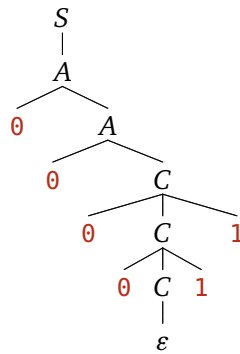
It is often useful to visualize derivations of strings in $L(G)$ using a *parse tree*. The parse tree for a string $w \in L(G)$ is a rooted ordered tree where

- Each leaf is labeled with a terminal or the empty string ϵ . Concatenating these in order from left to right yields the string w .
- Each internal node is labeled with a non-terminal. In particular, the root is labeled with the start non-terminal S .
- For each internal node v , there is a production rule $A \rightarrow \omega$ where A is the label of v and the symbols in ω are the labels of the children of v in order from left to right.

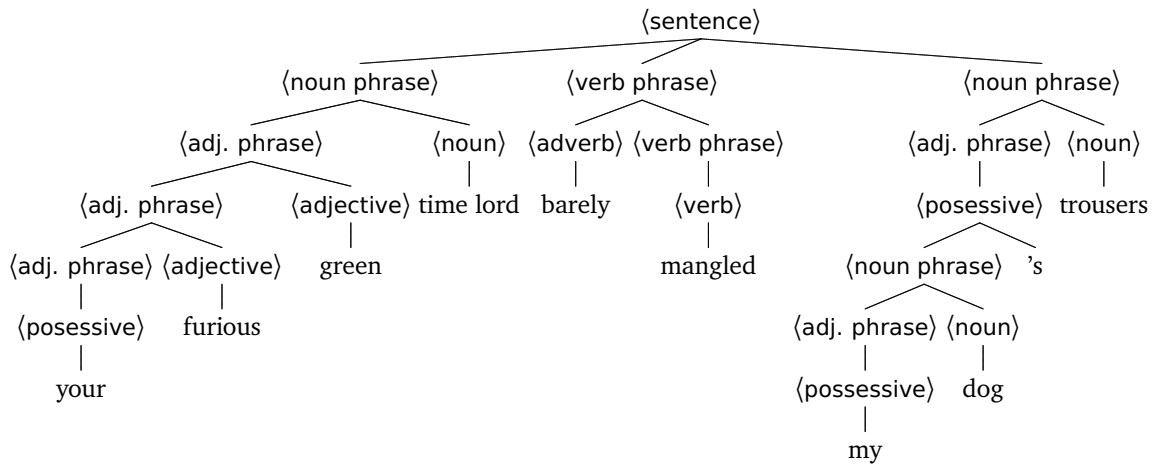
In other words, the production rules of the grammar describe *template trees* that can be assembled into larger parse trees. For example, the simple grammar on the previous page has the following templates, one for each production rule:



The same grammar gives us the following parse tree for the string 000011 :



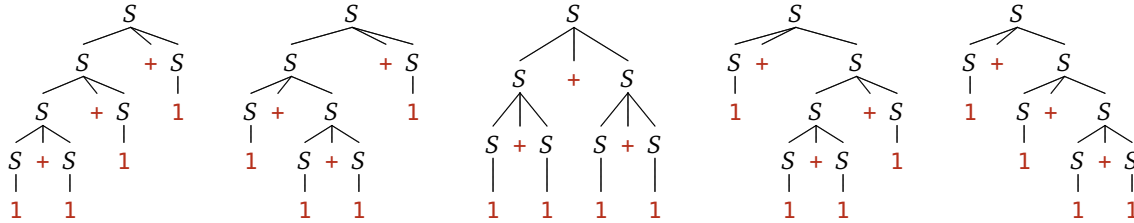
Our more complicated “English” grammar gives us parse trees like the following:



Any parse tree that contains at least one node with more than one non-terminal child corresponds to several different derivations. For example, when deriving an “English” sentence,

we have a choice of whether to expand the first ⟨noun phrase⟩ (“your furious green time lord”) before or after the second (“my dog’s trousers”).

A string w is **ambiguous** with respect to a grammar if there is more than one parse tree for w , and a grammar G is **ambiguous** if some string is ambiguous with respect to G . Neither of the previous example grammars is ambiguous. However, the grammar $S \rightarrow 1 \mid S+S$ is ambiguous, because the string $1+1+1+1$ has five different parse trees:



A context-free language L is **inherently ambiguous** if every context-free grammar that generates L is ambiguous. The language generated by the previous grammar (the regular language $(1+)^*1$) is *not* inherently ambiguous, because the unambiguous grammar $S \rightarrow 1 \mid 1+S$ generates the same language.

5.3 From Grammar to Language

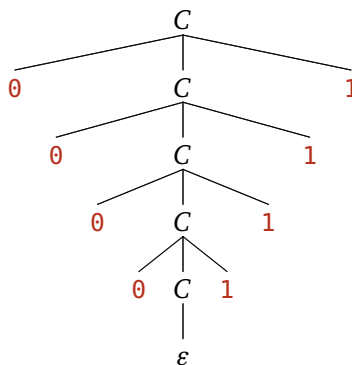
Let’s figure out the language generated by our first example grammar

$$S \rightarrow A \mid B \quad A \rightarrow 0A \mid 0C \quad B \rightarrow B1 \mid C1 \quad C \rightarrow \epsilon \mid 0C1.$$

Since the production rules for non-terminal C do not refer to any other non-terminal, let’s begin by figuring out $L(C)$. After playing around with the smaller grammar $C \rightarrow \epsilon \mid 0C1$ for a few seconds, you can probably guess that its language is $\{\epsilon, 01, 0011, 000111, \dots\}$, that is, the set all of strings of the form $0^n 1^n$ for some integer n . For example, we can derive the string 00001111 from the start non-terminal S using the following derivation:

$$C \rightsquigarrow 0C1 \rightsquigarrow 00C11 \rightsquigarrow 000C111 \rightsquigarrow 0000C1111 \rightsquigarrow 0000\epsilon 1111 = 00001111$$

The same derivation can be viewed as the following parse tree:



In fact, it is not hard to *prove* by induction that $L(C) = \{0^n 1^n \mid n \geq 0\}$ as follows. As usual when we prove that two sets X and Y are equal, the proof has two stages: one stage to prove $X \subseteq Y$, the other to prove $Y \subseteq X$.

- First we prove that $C \rightsquigarrow^* \theta^n 1^n$ for every non-negative integer n .

Fix an arbitrary non-negative integer n . Assume that $C \rightsquigarrow^* \theta^k 1^k$ for every non-negative integer $k < n$. There are two cases to consider.

- If $n = 0$, then $\theta^n 1^n = \varepsilon$. The rule $C \rightarrow \varepsilon$ implies that $C \rightsquigarrow \varepsilon$ and therefore $C \rightsquigarrow^* \varepsilon$.
- Suppose $n > 0$. The inductive hypothesis implies that $C \rightsquigarrow^* \theta^{n-1} 1^{n-1}$. Thus, the rule $C \rightarrow \theta C 1$ implies that $C \rightsquigarrow \theta C 1 \rightsquigarrow^* \theta(\theta^{n-1} 1^{n-1})1 = \theta^n 1^n$.

In both cases, we conclude that that $C \rightsquigarrow^* \theta^n 1^n$, as claimed.

- Next we prove that for every string $w \in \Sigma^*$ such that $C \rightsquigarrow^* w$, we have $w = \theta^n 1^n$ for some non-negative integer n .

Fix an arbitrary string w such that $C \rightsquigarrow^* w$. Assume that for any string x such that $|x| < |w|$ and $C \rightsquigarrow^* x$, we have $x = \theta^k 1^k$ for some non-negative integer k . There are two cases to consider, one for each production rule.

- If $w = \varepsilon$, then $w = \theta^0 1^0$.
- Suppose $w = \theta x 1$ for some string x such that $C \rightsquigarrow^* x$. Because $|x| = |w| - 2 < |w|$, the inductive hypothesis implies that $x = \theta^k 1^k$ for some integer k . Then we have $w = \theta^{k+1} 1^{k+1}$.

In both cases, we conclude that that $w = \theta^n 1^n$ for some non-negative integer n , as claimed.

The first proof uses induction on strings, following the boilerplate proposed in the previous lecture; in particular, the case analysis mirrors the recursive definition of “string”. The second proof uses *structural induction* on the grammar; the case analysis mirrors the recursive definition of the language of S , as described by the production rules. In both proofs, the inductive hypothesis is “Assume there is no smaller counterexample.”

Similar analysis implies that $L(A) = \{\theta^m 1^n \mid m > n\}$ and $L(B) = \{\theta^m 1^n \mid m < n\}$, and therefore $L(S) = \{\theta^m 1^n \mid m \neq n\}$.

5.4 More Examples



Give three or four examples of simple but interesting context-free grammars. Some possibilities:

- Same number of 0s and 1s
- Different number of 0s and 1s
- Palindromes
- Balanced parentheses
- Arithmetic/algebraic expressions
- Regular expressions

5.5 Regular Languages are Context-Free

The following inductive argument proves that every regular language is also a context-free language. Let L be an arbitrary regular language, encoded by some regular expression R . Assume that any regular expression shorter than R represents a context-free language. (“Assume no smaller counterexample.”) We construct a context-free grammar for L as follows. There are several cases to consider.

- Suppose L is empty. Then L is generated by the trivial grammar $S \rightarrow S$.
- Suppose $L = \{w\}$ for some string $w \in \Sigma^*$. Then L is generated by the grammar $S \rightarrow w$.
- Suppose L is the union of some regular languages L_1 and L_2 . The inductive hypothesis implies that L_1 and L_2 are context-free. Let G_1 be a context-free language for L_1 with starting non-terminal S_1 , and let G_2 be a context-free language for L_2 with starting non-terminal S_2 , where the non-terminal sets in G_1 and G_2 are disjoint. Then $L = L_1 \cup L_2$ is generated by the production rule $S \rightarrow S_1 \mid S_2$.
- Suppose L is the concatenation of some regular languages L_1 and L_2 . The inductive hypothesis implies that L_1 and L_2 are context-free. Let G_1 be a context-free language for L_1 with starting non-terminal S_1 , and let G_2 be a context-free language for L_2 with starting non-terminal S_2 , where the non-terminal sets in G_1 and G_2 are disjoint. Then $L = L_1 L_2$ is generated by the production rule $S \rightarrow S_1 S_2$.
- Suppose L is the Kleene closure of some regular language L_1 . The inductive hypothesis implies that L_1 is context-free. Let G_1 be a context-free language for L_1 with starting non-terminal S_1 . Then $L = L_1^*$ is generated by the production rule $S \rightarrow \varepsilon \mid S_1 S$.

In every case, we have found a context-free grammar that generates L , which means L is context-free.

In the next lecture note, we will prove that the context-free language $\{0^n 1^n \mid n \geq 0\}$ is not regular. (In fact, this is the *canonical example* of a non-regular language.) Thus, context-free grammars are strictly more expressive than regular expressions.

5.6 Not Every Language is Context-Free

Again, you may be tempted to conjecture that *every* language is context-free, but a variant of our earlier cardinality argument implies that this is not the case.

Any context-free grammar over the alphabet Σ can be encoded as a string over the alphabet $\Sigma \cup \Gamma \cup \{\varepsilon, \rightarrow, |, \$\}$, where $\$$ indicates the end of the production rules for each non-terminal. For example, our example grammar

$$S \rightarrow A \mid B \quad A \rightarrow 0A \mid 0C \quad B \rightarrow B1 \mid C1 \quad C \rightarrow \varepsilon \mid 0C1$$

can be encoded as the string

$$S \rightarrow A \mid B \$ A \rightarrow 0A \mid 0C \$ B \rightarrow B1 \mid C1 \$ C \rightarrow \varepsilon \mid 0C1 \$$$

We can further encode any such string as a *binary* string by associating each symbol in the set $\Sigma \cup \Gamma \cup \{\varepsilon, \rightarrow, |, \$\}$ with a different binary substring. Specifically, if we encode each of the grammar symbols $\varepsilon, \rightarrow, |, \$$ as a string of the form 11^*0 , each terminal in Σ as a string of the form 011^*0 , and each non-terminal as a string of the form 0011^*0 , we can unambiguously recover the grammar from the encoding. For example, applying the code

$$\begin{array}{lll} \varepsilon \mapsto 10 & 0 \mapsto 010 & S \mapsto 0010 \\ \rightarrow \mapsto 110 & 1 \mapsto 0110 & A \mapsto 00110 \\ | \mapsto 1110 & & B \mapsto 001110 \\ \$ \mapsto 11110 & & C \mapsto 0011110 \end{array}$$

transforms our example grammar into the 135-bit string

```

00101100011011100011101111000110
11001000110111001000111101111000
11101100011100110111000111100110
11110001111011010111001000111100
1011110.

```

Adding a **1** to the start of this bit string gives us the binary encoding of the integer

51 115 617 766 581 763 757 672 062 401 233 529 937 502.

Our construction guarantees that two different context-free grammars over the same language (ignoring changing the names of the non-terminals) yield different positive integers. Thus, the set of context-free grammars over any alphabet is *at most* as large as the set of integers, and is therefore countably infinite. (Most integers are not encodings of context-free grammars, but that only helps us.) It follows that the set of all context-free *languages* over any fixed alphabet is also countably infinite. But we already showed that the set of *all* languages over any alphabet is uncountably infinite. So almost all languages are non-context-free!

Although we will probably not see them in this course, there are techniques for proving that certain languages are not context-free, just as there are for proving certain languages are not regular. In particular, the $\{0^n 1^n 0^n \mid n \geq 0\}$ is not context-free. (In fact, this is the *canonical example* of a non-context-free language.)

*5.7 Recursive Automata

All the flavors of finite-state automata we have seen so far describe/encode/accept/compute *regular* languages; these are precisely the languages that can be constructed from individual strings by union, concatenation, and unbounded repetition. Just as context-free grammars are recursive generalizations of regular expressions, we can define a class of machines called *recursive automata*, which generalize (nondeterministic) finite-state automata.

Formally, a **recursive automaton** consists of the following components:

- A non-empty finite set Σ , called the **input alphabet**
- Another non-empty finite set N , disjoint from Σ , whose elements are called **module names**
- A **start name** $S \in N$
- A set $M = \{M_A \mid A \in N\}$ of NFAs over the alphabet $\Sigma \cup N$ called **modules**, each with a single accepting state. Each module M_A has the following components:
 - A finite set Q_A of **states**, such that $Q_A \cap Q_B = \emptyset$ for all $A \neq B$
 - A **start** state $s_A \in Q_A$
 - A **terminal** or **accepting** state $t_A \in Q_A$
 - A **transition function** $\delta_A: Q_A \times (\Sigma \cup \{\varepsilon\} \cup N) \rightarrow 2^{Q_A}$.

Equivalently, we have a single global transition function $\delta: Q \times (\Sigma \cup \{\varepsilon\} \cup N) \rightarrow 2^Q$, where $Q = \bigcup_{A \in N} Q_A$, such that for any name A and any state $q \in Q_A$ we have $\delta(q) \subseteq Q_A$. Machine M_S is called the **main module**.

A **configuration** of a recursive automaton is a triple (w, q, s) , where w is a string in Σ^* called the **input**, q is a state in Q called the **local state**, and s is a string in Q^* called the **stack**. The module containing the local state q is called the **active module**. A configuration can be changed by three types of transitions.

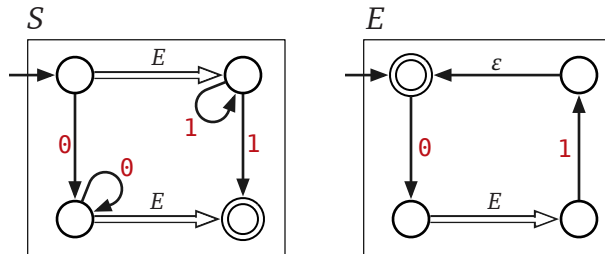
- A **read** transition consumes the first symbol in the input and changes the local state within the current module, just like a standard NFA.
- An **epsilon** transition changes the local state within the current module, without consuming any input characters, just like a standard NFA.
- A **call** transition chooses an arbitrary name A , changes the current state q_0 to some state in $\delta(q, A)$, and pushes the corresponding start state s_A onto the stack (thereby changing the active module to M_A), without consuming any input characters.
- Finally, if the current state is the terminal state of some module *and* the stack is non-empty, a **return** transition pops the top state off the stack and makes it the new local state (thereby possibly changing the active module), without consuming any input characters.

Symbolically, we can describe these transitions as follows:

$$\begin{array}{lll}
 \text{read:} & (ax, q, \sigma) \mapsto (x, q', \sigma) & \text{for some } q' \in \delta(q, a) \\
 \text{epsilon:} & (w, q, \sigma) \mapsto (w, q', \sigma) & \text{for some } q' \in \delta(q, \varepsilon) \\
 \text{call:} & (w, q, \sigma) \mapsto (w, s_A, q' \cdot \sigma) & \text{for some } A \in N \text{ and } q' \in \delta(q, A) \\
 \text{return:} & (w, t_A, q \cdot \sigma) \mapsto (w, q, \sigma) &
 \end{array}$$

A recursive automaton **accepts** a string w if there is a *finite* sequence of transitions starting at the start configuration (w, s_S, ε) and ending at the terminal configuration $(\varepsilon, t_S, \varepsilon)$.

For example, the following recursive automaton accepts the language $\{0^m 1^n \mid m \neq n\}$. The recursive automaton has two component machines; the start machine named S and a “subroutine” named E (for “equal”) that accepts the language $\{0^n 1^n \mid n \geq 0\}$. White arrows indicate recursive transitions.



A recursive automaton for the language $\{0^m 1^n \mid m \neq n\}$

Lemma 5.1. *Every context-free language is accepted by a recursive automaton.*

Proof:



Direct construction from the CFG, with one module per nonterminal.

□

For example, the context-free grammar

$$\begin{array}{l}
 S \rightarrow 0A \mid B1 \\
 A \rightarrow 0A \mid E \\
 B \rightarrow B1 \mid E \\
 E \rightarrow \varepsilon \mid 0E0
 \end{array}$$

leads to the following recursive automaton with four modules:



Figure!

Lemma 5.2. *Every recursive automaton accepts a context-free language.*

Proof (sketch): Let $R = (\Sigma, N, S, \delta, M)$ be an arbitrary recursive automaton. We define a context-free grammar G that describes the language accepted by R as follows.

The set of nonterminals in G is isomorphic to the state set Q ; that is, for each state $q \in Q$, the grammar contains a corresponding nonterminal $[q]$. The language of $[q]$ will be the set of strings w such that there is a finite sequence of transitions starting at the start configuration (w, q, ε) and ending at the terminal configuration $(\varepsilon, t, \varepsilon)$, where t is the terminal state of the module containing q .

The grammar has four types of production rules, corresponding to the four types of transitions:

- **read:** For each symbol a and each pair of states p and q such that $p \in \delta(q, a)$, the grammar contains the production rule $[q] \rightarrow a[p]$.
- **epsilon:** For any two states p and q such that $p \in \delta(q, \varepsilon)$, the grammar contains the production rule $[q] \rightarrow [p]$.
- **call:** Each name A and each pair of states p and q such that $p \in \delta(q, A)$, the grammar contains the production rule $[q] \rightarrow [s_A][p]$.
- **return:** Each name A , the grammar contains the production rule $[t_A] \rightarrow \varepsilon$.

Finally, the starting nonterminal of G is $[s_S]$, which corresponds to the start state of the main module.

We can now argue inductively that the grammar G and the recursive automaton R describe the same language. Specifically, any sequence of transitions in R from (w, s_S, ε) to $(\varepsilon, t_S, \varepsilon)$ can be transformed mechanically into a derivation of w from the nonterminal $[s_S]$ in G . Symmetrically, the **leftmost** derivation of any string w in G can be mechanically transformed into an accepting sequence of transitions in R . We omit the straightforward but tedious details. \square

For example, the recursive automaton on the previous page gives us the following context-free grammar. To make the grammar more readable, I've renamed the nonterminals corresponding to start and terminal states: $S = [s_S]$, $T = [t_S]$, and $E = [s_E] = [t_E]$:

$$\begin{array}{ll}
 S \rightarrow EA \mid \emptyset B & E \rightarrow \varepsilon \mid \emptyset X \\
 A \rightarrow \mathbf{1}A \mid \mathbf{1}T & X \rightarrow EY \\
 B \rightarrow \emptyset B \mid ET & Y \rightarrow \mathbf{1}Z \\
 T \rightarrow \varepsilon & Z \rightarrow E
 \end{array}$$

Our earlier proofs imply that we can forbid ε -transitions or even allow regular-expression transitions in our recursive automata without changing the set of languages they accept.

*5.8 Chomsky Normal Form

For many algorithmic problems involving context-free grammars, it is helpful to consider grammars with a particular special structure called **Chomsky normal form**, abbreviated **CNF**:

- The starting non-terminal S does not appear on the right side of any production rule.

- The starting non-terminal S may have the production rule $S \rightarrow \varepsilon$.
- The right side of every other production rule is either a single terminal symbol or a string of exactly two non-terminals—that is, every other production rule has the form $A \rightarrow BC$ or $A \rightarrow a$.

A particularly attractive feature of CNF grammars is that they yield *full binary* parse trees; in particular, every parse tree for a string of length $n > 0$ has exactly $2n - 1$ non-terminal nodes. Consequently, any string of length n in the language of a CNF grammar can be derived in exactly $2n - 1$ production steps. It follows that we can actually determine whether a string belongs to the language of a CNF grammar by brute-force consideration of all possible derivations of the appropriate length.

For arbitrary context-free grammars, there is no similar upper bound on the length of a derivation, and therefore no similar brute-force membership algorithm, because the grammar may contain additional *ε -productions* of the form $A \rightarrow \varepsilon$ and/or *unit productions* of the form $A \rightarrow B$, where both A and B are non-terminals. Unit productions introduce nodes of degree 1 into any parse tree, and ε -productions introduce leaves that do not contribute to the word being parsed.

Fortunately, it is possible to determine membership in the language of an arbitrary context-free grammar, thanks to the following theorem. Two context-free grammars are *equivalent* if they define the same language.

Every context-free grammar is equivalent to a grammar in Chomsky normal form.

To be more specific, define the *total length* of a context-free grammar to be the number of symbols needed to write down the grammar; up to constant factors, the total length is the sum of the lengths of the production rules.

Theorem 5.3. *For every context-free grammar with total length L , there is an equivalent grammar in Chomsky normal form with total length $O(L^2)$, which can be computed in $O(L^2)$ time.*

Converting an arbitrary grammar into Chomsky normal form is a complex task. Fortunately, for most applications of context-free grammars, it's enough to know that the algorithm exists. For the sake of completeness, however, I will describe one such conversion algorithm here. This algorithm consists of several relatively straightforward stages. Efficient implementation of some of these stages requires standard graph-traversal algorithms, which we will describe much later in the course.

0. Add a new starting non-terminal. Add a new non-terminal S' and a production rule $S' \rightarrow S$, where S is the starting non-terminal for the given grammar. S' will be the starting non-terminal for the resulting CNF grammar. (In fact, this step is necessary only when $S \rightsquigarrow^* \varepsilon$, but at this point in the conversion process, we don't yet know whether that's true.)

1. Decompose long production rules. For each production rule $A \rightarrow \omega$ whose right side w has length greater than two, add new production rules of length two that still permit the derivation $A \rightsquigarrow^* \omega$. Specifically, suppose $\omega = \alpha\chi$ for some symbol $\alpha \in \Sigma \cup \Gamma$ and string $\chi \in (\Sigma \cup \Gamma)^*$. The algorithm replaces $A \rightarrow \omega$ with two new production rules $A \rightarrow \alpha B$ and $B \rightarrow \chi$, where B is a new non-terminal, and then (if necessary) recursively decomposes the production rule $B \rightarrow \chi$. For

example, we would replace the long production rule $A \rightarrow 0BC1CB$ with the following sequence of short production rules, where each A_i is a new non-terminal:

$$A \rightarrow 0A_1 \quad A_1 \rightarrow BA_2 \quad A_2 \rightarrow CA_3 \quad A_3 \rightarrow 1A_4 \quad A_4 \rightarrow CB$$

This stage can significantly increase the number of non-terminals and production rules, but it increases the *total length* of all production rules by at most a small constant factor.² Moreover, for the remainder of the conversion algorithm, every production rule has length at most two. The running time of this stage is $O(L)$.

2. Identify nullable non-terminals. A non-terminal A is *nullable* if and only if $A \rightsquigarrow^* \varepsilon$. The recursive definition of \rightsquigarrow^* implies that A is nullable if and only if the grammar contains a production rule $A \rightarrow \omega$ where ω consists entirely of nullable non-terminals (in particular, if $\omega = \varepsilon$). You may be tempted to transform this recursive characterization directly into a recursive algorithm, but this is a bad idea; the resulting algorithm would fall into an infinite loop if (for example) the same non-terminal appeared on both sides of the same production rule. Instead, we apply the following *fixed-point* algorithm, which repeatedly scans through the entire grammar until a complete scan discovers no new nullable non-terminals.

```

NULLABLES( $\Sigma, \Gamma, R, S$ ):
 $\Gamma_\varepsilon \leftarrow \emptyset$      $\langle\langle$ known nullable non-terminals $\rangle\rangle$ 
 $done \leftarrow \text{FALSE}$ 
while  $\neg done$ 
     $done \leftarrow \text{TRUE}$ 
    for each non-terminal  $A \in \Gamma \setminus \Gamma_\varepsilon$ 
        for each production rule  $A \rightarrow \omega$ 
            if  $\omega \in \Gamma_\varepsilon^*$ 
                add  $A$  to  $\Gamma_\varepsilon$ 
                 $done \leftarrow \text{FALSE}$ 
return  $\Gamma_\varepsilon$ 

```

At this point in the conversion algorithm, if S' is *not* identified as nullable, then we can safely remove it from the grammar and use the original starting nonterminal S instead.

As written, NULLABLES runs in $O(nL) = O(L^2)$ time, where n is the number of non-terminals in Γ . Each iteration of the main loop except the last adds at least one non-terminal to Γ_ε , so the algorithm halts after at most $n + 1 \leq L$ iterations, and in each iteration, we examine at most L production rules. There is a faster implementation of NULLABLES that runs in $O(n + L) = O(L)$ time,³ but since other parts of the conversion algorithm already require $O(L^2)$ time, we needn't bother.

3. Eliminate ε -productions. First, remove every production rule of the form $A \rightarrow \varepsilon$. Then for each production rule $A \rightarrow w$, add all possible new production rules of the form $A \rightarrow w'$, where w'

²In most textbook descriptions of this conversion algorithm, this stage is performed *last*, after removing ε -productions and unit productions. But with the stages in that traditional order, removing ε -productions could *exponentially* increase the length of the grammar in the worst case! Consider the production rule $A \rightarrow (BC)^k$, where B is nullable but C is not. Decomposing this rule first and then removing ε -productions introduces about $3k$ new rules; whereas, removing ε -productions first introduces 2^k new rules, most of which then must then be further decomposed.

³Consider the bipartite graph whose vertices correspond to non-terminals and the right sides of production rules, with one edge per rule. The faster algorithm is a modified breadth-first search of this graph, starting at the vertex representing ε .

is a **non-empty** string obtained from w by removing one nullable non-terminal. For example, if the grammar contained the production rule $A \rightarrow BC$, where B and C are both nullable, we would add two new production rules $A \rightarrow B \mid C$. Finally, if the starting nonterminal S' was identified as nullable in the previous stage, add the production rule $S' \rightarrow \varepsilon$; this will be the *only* ε -production in the final grammar. This phase of the conversion runs in $O(L)$ time and at most triples the number of production rules.

4. Merge equivalent non-terminals. We say that two non-terminals A and B are **equivalent** if they can be derived from each other: $A \rightsquigarrow^* B$ and $B \rightsquigarrow^* A$. Because we have already removed ε -productions, any such derivation must consist entirely of unit productions. For example, in the grammar

$$S \rightarrow B \mid C, \quad A \rightarrow B \mid D \mid CC \mid \emptyset, \quad B \rightarrow C \mid AD \mid \mathbf{1}, \quad C \rightarrow A \mid DA, \quad D \rightarrow BA \mid CS,$$

non-terminals A, B, C are all equivalent, but S is not in that equivalence class (because we cannot derive S from A) and neither is D (because we cannot derive A from D).

Construct a directed graph G whose vertices are the non-terminals and whose edges correspond to unit productions, in $O(L)$ time. Then two non-terminals are equivalent if and only if they are in the same strong component of G . Compute the strong components of G in $O(L)$ time using, for example, the algorithm of Kosaraju and Sharir. Then merge all the non-terminals in each equivalence class into a single non-terminal. Finally, remove any unit productions of the form $A \rightarrow A$. The total running time for this phase is $O(L)$. Starting with our example grammar above, merging B and C with A and removing the production $A \rightarrow A$ gives us the simpler grammar

$$S \rightarrow A, \quad A \rightarrow AA \mid D \mid DA \mid \emptyset \mid \mathbf{1}, \quad D \rightarrow AA \mid AS.$$

We could further simplify the grammar by merging all non-terminals reachable from S using only unit productions (in this case, merging non-terminals S and S), but this further simplification is unnecessary.

5. Remove unit productions. Once again, we construct a directed graph G whose vertices are the non-terminals and whose edges correspond to unit productions, in $O(L)$ time. Because no two non-terminals are equivalent, G is acyclic. Thus, using topological sort, we can index the non-terminals A_1, A_2, \dots, A_n such that for every unit production $A_i \rightarrow A_j$ we have $i < j$, again in $O(L)$ time; moreover, we can assume that the starting non-terminal is A_1 . (In fact, both the dag G and the linear ordering of non-terminals was already computed in the previous phase!!)

Then for each index j in decreasing order, for each unit production $A_i \rightarrow A_j$ and each production $A_j \rightarrow \omega$, we add a new production rule $A_i \rightarrow \omega$. At this point, all unit productions are redundant and can be removed. Applying this algorithm to our example grammar above gives us the grammar

$$S \rightarrow AA \mid AS \mid DA \mid \emptyset \mid \mathbf{1}, \quad A \rightarrow AA \mid AS \mid DA \mid \emptyset \mid \mathbf{1}, \quad D \rightarrow AA \mid AS.$$

In the worst case, each production rule for A_n is copied to each of the other $n - 1$ non-terminals. Thus, this phase runs in $\Theta(nL) = O(L^2)$ time and increases the length of the grammar to $\Theta(nL) = O(L^2)$ in the worst case.

This phase dominates the running time of the CNF conversion algorithm. Unlike previous phases, no faster algorithm for removing unit transformations is known! There are grammars of length L with unit productions such that any equivalent grammar without unit productions has

length $\Omega(L^{1.499999})$ (for any desired number of 9s), but this lower bound does not rule out the possibility of an algorithm that runs in only $O(L^{3/2})$ time. Closing the gap between $\Omega(L^{3/2-\epsilon})$ and $O(L^2)$ has been an open problem since the early 1980s.

6. Protect terminals. Finally, for each terminal $a \in \Sigma$, we introduce a new non-terminal A_a and a new production rule $A_a \rightarrow a$, and then replace a with A_a in every production rule of length 2. This completes the conversion to Chomsky normal form. As claimed, the total running time of the algorithm is $O(L^2)$, and the total length of the output grammar is also $O(L^2)$.

CNF Conversion Example

As a running example, let's apply these stages one at a time to our first example grammar.

$$S \rightarrow A | B \quad A \rightarrow \emptyset A | \emptyset C \quad B \rightarrow B1 | C1 \quad C \rightarrow \varepsilon | \emptyset C1$$

- o. Add a new starting non-terminal S' .

$$\underline{S'} \rightarrow S \quad S \rightarrow A | B \quad A \rightarrow \emptyset A | \emptyset C \quad B \rightarrow B1 | C1 \quad C \rightarrow \varepsilon | \emptyset C1$$

1. Decompose the long production rule $C \rightarrow \emptyset C1$.

$$S' \rightarrow S \quad S \rightarrow A | B \quad A \rightarrow \emptyset A | \emptyset C \quad B \rightarrow B1 | C1 \quad \underline{C \rightarrow \varepsilon | \emptyset D} \quad \underline{D \rightarrow C1}$$

2. Identify C as the only nullable non-terminal. Because S' is not nullable, remove the production rule $S' \rightarrow S$.
3. Eliminate the ε -production $C \rightarrow \varepsilon$.

$$S \rightarrow A | B \quad A \rightarrow \emptyset A | \emptyset C | \underline{\emptyset} \quad B \rightarrow B1 | C1 | \underline{1} \quad C \rightarrow \emptyset D \quad D \rightarrow C1 | \underline{1}$$

4. No two non-terminals are equivalent, so there's nothing to merge.
5. Remove the unit productions $S' \rightarrow S$, $S \rightarrow A$, and $S \rightarrow B$.

$$S \rightarrow \underline{\emptyset A | \emptyset C | B1 | C1 | \emptyset | 1}$$

$$A \rightarrow \emptyset A | \emptyset C | \emptyset \quad B \rightarrow B1 | C1 | 1 \quad C \rightarrow \emptyset D \quad D \rightarrow C1 | 1.$$

6. Finally, protect the terminals \emptyset and 1 to obtain the final CNF grammar.

$$\underline{S \rightarrow EA | EC | BF | CF | \emptyset | 1}$$

$$A \rightarrow \underline{EA | EC} | \emptyset \quad B \rightarrow \underline{BF | CF} | 1$$

$$C \rightarrow \underline{ED} \quad D \rightarrow \underline{CF} | 1$$

$$\underline{E \rightarrow \emptyset} \quad \underline{F \rightarrow 1}$$

Exercises

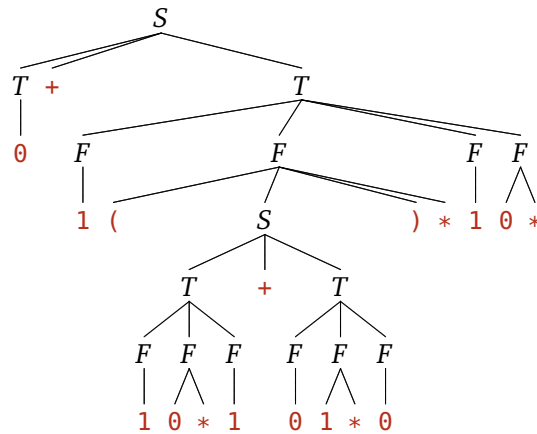
1. Describe context-free grammars that generate each of the following languages. The function $\#(x, w)$ returns the number of occurrences of the **substring** x in the string w . For example, $\#(0, 101001) = 3$ and $\#(010, 1010100011) = 2$.
 - (a) All strings in $\{0, 1\}^*$ whose length is divisible by 5.
 - (b) All strings in $\{0, 1\}^*$ representing a non-negative multiple of 5 in binary.
 - (c) $\{w \in \{0, 1\}^* \mid \#(0, w) = \#(1, w)\}$
 - (d) $\{w \in \{0, 1\}^* \mid \#(0, w) \neq \#(1, w)\}$
 - (e) $\{w \in \{0, 1\}^* \mid \#(00, w) = \#(11, w)\}$
 - (f) $\{w \in \{0, 1\}^* \mid \#(01, w) = \#(10, w)\}$
 - (g) $\{w \in \{0, 1\}^* \mid \#(0, w) = \#(1, w) \text{ and } |w| \text{ is a multiple of } 3\}$
 - (h) $\{0, 1\}^* \setminus \{0^n 1^n \mid n \geq 0\}$
 - (i) $\{0^n 1^{2n} \mid n \geq 0\}$
 - (j) $\{0, 1\}^* \setminus \{0^n 1^{2n} \mid n \geq 0\}$
 - (k) $\{0^n 1^m \mid 0 \leq 2m \leq n < 3m\}$
 - (l) $\{0^i 1^j 2^{i+j} \mid i, j \geq 0\}$
 - (m) $\{0^i 1^j 2^k \mid i = j \text{ or } j = k\}$
 - (n) $\{0^i 1^j 2^k \mid i \neq j \text{ or } j \neq k\}$
 - (o) $\{0^i 1^j 0^j 1^i \mid i, j \geq 0\}$
 - (p) $\{w \$ 0^{\#(0, w)} \mid w \in \{0, 1\}^*\}$
 - (q) $\{xy \mid x, y \in \{0, 1\}^* \text{ and } x \neq y \text{ and } |x| = |y|\}$
 - (r) $\{x \$ y^R \mid x, y \in \{0, 1\}^* \text{ and } x \neq y\}$
 - (s) $\{x \$ y \mid x, y \in \{0, 1\}^* \text{ and } \#(0, x) = \#(1, y)\}$
 - (t) $\{0, 1\}^* \setminus \{ww \mid w \in \{0, 1\}^*\}$
 - (u) All strings in $\{0, 1\}^*$ that are *not* palindromes.
 - (v) All strings in $\{(,), \diamond\}^*$ in which the parentheses are balanced and the symbol \diamond appears at most four times. For example, $()(())$ and $(\diamond((())\diamond)(())\diamond)$ and $\diamond\diamond$ are strings in this language, but $)((()$ and $(\diamond\diamond)\diamond$ are not.
2. Describe recursive automata for each of the languages in problem 1. (“Describe” does not necessarily mean “draw”!)
3. Prove that if L is a context-free language, then L^R is also a context-free language. [Hint: How do you reverse a context-free grammar?]
4. Consider a generalization of context-free grammars that allows any *regular expression* over $\Sigma \cup \Gamma$ to appear on the right side of a production rule. Without loss of generality, for each non-terminal $A \in \Gamma$, the generalized grammar contains a single regular expression $R(A)$. To apply a production rule to a string, we replace any non-terminal A with an arbitrary word

in the language described by $R(A)$. As usual, the language of the generalized grammar is the set of all strings that can be derived from its start non-terminal.

For example:, the following generalized context-free grammar describes the language of all regular expressions over the alphabet $\{0, 1\}$:

$$\begin{aligned} S &\rightarrow (T+)^*T + \emptyset && \text{(Regular expressions)} \\ T &\rightarrow \varepsilon + F^*F && \text{(Terms = summable expressions)} \\ F &\rightarrow (0 + 1 + (S))^*(\varepsilon) && \text{(Factors = concatenable expressions)} \end{aligned}$$

Here is a parse tree for the regular expression $0+1(10^*1+01^*0)^*10^*$ (which represents the set of all binary numbers divisible by 3):



Prove that every *generalized* context-free grammar describes a context-free language. In other words, show that allowing regular expressions to appear in production rules does not increase the expressive power of context-free grammars.