

9

Restricted computational models

“Happy families are all alike; every unhappy family is unhappy in its own way”, Leo Tolstoy (opening of the book “Anna Karenina”).

We have seen that a great many models of computation are *Turing equivalent*, including our NAND++/NAND« programs and Turing machines, standard programming languages such as C/Python/-Javascript etc., and other models such as the λ calculus and even the game of life. The flip side of this is that for all these models, Rice’s theorem (see [Section 8.4.1](#)) holds as well, which means that deciding any semantic property of programs in such a model is *uncomputable*.

The uncomputability of halting and other semantic specification problems for Turing equivalent models motivates coming up with **restricted computational models** that are **(a)** powerful enough to capture a set of functions useful for certain applications but **(b)** weak enough that we can still solve semantic specification problems on them. In this chapter we will discuss several such examples.

9.1 TURING COMPLETENESS AS A BUG

We have seen that seemingly simple computational models or systems can turn out to be Turing complete. The [following webpage](#) lists several examples of formalisms that “accidentally” turned out to be Turing complete, including supposedly limited languages such as the C preprocessor, CCS, SQL, sendmail configuration, as well as games such as Minecraft, Super Mario, and the card game “Magic: The gathering”. This is not always a good thing, as it means that such formalisms can give rise to arbitrarily complex behavior. For example, the postscript format (a precursor of PDF) is a Turing-complete programming language meant to describe documents for printing. The expressive power of postscript can allow for short description of very complex images. But it also gives rise to some nasty surprises, such as

Learning Objectives:

- See that Turing completeness is not always a good thing
- Two important examples of non-Turing-complete, always-halting formalisms: *regular expressions* and *context-free grammars*.
- The pumping lemmas for both these formalisms, and examples of non regular and non context-free functions.
- Examples of computable and uncomputable *semantic properties* of regular expressions and context-free grammars.

the attacks described in [this page](#) ranging from using infinite loops as a denial of service attack, to accessing the printer's file system.

An interesting recent example of the pitfalls of Turing-completeness arose in the context of the cryptocurrency [Ethereum](#). The distinguishing feature of this currency is the ability to design "smart contracts" using an expressive (and in particular Turing-complete) language. In our current "human operated" economy, Alice and Bob might sign a contract to agree that if condition X happens then they will jointly invest in Charlie's company. Ethereum allows Alice and Bob to create a joint venture where Alice and Bob pool their funds together into an account that will be governed by some program P that decides under what conditions it disburses funds from it. For example, one could imagine a piece of code that interacts between Alice, Bob, and some program running on Bob's car that allows Alice to rent out Bob's car without any human intervention or overhead.

Specifically Ethereum uses the Turing-complete programming [solidity](#) which has a syntax similar to Javascript. The flagship of Ethereum was an experiment known as The "Decentralized Autonomous Organization" or [The DAO](#). The idea was to create a smart contract that would create an autonomously run decentralized venture capital fund, without human managers, were shareholders could decide on investment opportunities. The DAO was the biggest crowdfunding success in history and at its height was worth 150 million dollars, which was more than ten percent of the total Ethereum market. Investing in the DAO (or entering any other "smart contract") amounts to providing your funds to be run by a computer program. i.e., "code is law", or to use the words the DAO described itself: "The DAO is borne from immutable, unstoppable, and irrefutable computer code". Unfortunately, it turns out that (as we saw in [Chapter 8](#)) understanding the behavior of Turing-complete computer programs is quite a hard thing to do. A hacker (or perhaps, some would say, a savvy investor) was able to fashion an input that would cause the DAO code to essentially enter into an infinite recursive loop in which it continuously transferred funds into their account, thereby [cleaning out about 60 million dollars](#) out of the DAO. While this transaction was "legal" in the sense that it complied with the code of the smart contract, it was obviously not what the humans who wrote this code had in mind. There was a lot of debate in the Ethereum community how to handle this, including some partially successful "Robin Hood" attempts to use the same loophole to drain the DAO funds into a secure account. Eventually it turned out that the code is mutable, stoppable, and refutable after all, and the Ethereum community decided to do a "hard fork" (also known as a "bailout") to revert history to before this transaction. Some elements of the community strongly opposed this

decision, and so an alternative currency called **Ethereum Classic** was created that preserved the original history.

9.2 REGULAR EXPRESSIONS

One of the most common tasks in computing is to *search* for a piece of text. At its heart, the *search problem* is quite simple. The user gives out a function $F : \{0, 1\}^* \rightarrow \{0, 1\}$, and the system applies this function to a set of candidates $\{x_0, \dots, x_k\}$, returning all the x_i 's such that $F(x_i) = 1$. However, we typically do not want the system to get into an infinite loop just trying to evaluate this function! For this reason, such systems often do not allow the user to specify an *arbitrary* function using some Turing-complete formalism, but rather a function that is described by a restricted computational model, and in particular one in which all functions halt. One of the most popular models for this application is the model of **regular expressions**. You have probably come across regular expressions if you ever used an advanced text editor, a command line shell, or have done any kind of manipulations of text files.

A *regular expression* over some alphabet Σ is obtained by combining elements of Σ with the operation of concatenation, as well as $|$ (corresponding to *or*) and $*$ (corresponding to repetition zero or more times).¹ For example, the following regular expression over the alphabet $\{0, 1\}$ corresponds to the set of all even length strings $x \in \{0, 1\}^*$ where the digit at location $2i$ is the same as the one at location $2i + 1$ for every i :

$$(00|11)^* \quad (9.1)$$

The following regular expression over the alphabet $\{a, b, c, d, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ corresponds to the set of all strings that consist of a sequence of one or more of the letters a - b followed by a sequence of one or more digits (without a leading zero):

$$(a|b|c|d)(a|b|c|d)^*(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^* \quad (9.2)$$

Formally, regular expressions are defined by the following recursive definition:²

Definition 9.1 — Regular expression. A *regular expression* exp over an alphabet Σ is a string over $\Sigma \cup \{(\, , |, *, \emptyset, ""\}$ that has one of the following forms:

1. $exp = \sigma$ where $\sigma \in \Sigma$

¹ Common implementations of regular expressions in programming languages and shells typically include some extra operations on top of $|$ and $*$, but these can all be implemented as “syntactic sugar” using the operators $|$ and $*$.

² We have seen a recursive definition before in the setting of λ expressions (Definition 7.4). Just like recursive functions, we can define a concept recursively. A definition of some class \mathcal{C} of objects can be thought of as defining a function that maps an object o to either *VALID* or *INVALID* depending on whether $o \in \mathcal{C}$. Thus we can think of Definition 9.1 as defining a recursive function that maps a string exp over $\Sigma \cup \{(\, , |, *, \emptyset, ""\}$ to *VALID* or *INVALID* depending on whether exp describes a valid regular expression.

2. $exp = (exp' | exp'')$ where exp', exp'' are regular expressions.
3. $exp = (exp')(exp'')$ where exp', exp'' are regular expressions. (We often drop the parenthesis when there is no danger of confusion and so write this as $exp exp'$.)
4. $exp = (exp')^*$ where exp' is a regular expression.

Finally we also allow the following “edge cases”: $exp = \emptyset$ and $exp = \cdot$.³

Every regular expression exp correspond to a function $\Phi_{exp} : \Sigma^* \rightarrow \{0, 1\}$ where $\Phi_{exp}(x) = 1$ if x matches the regular expression. The definition of “matching” is recursive as well. For example, if exp and exp' match the strings x and x' , then the expression $exp exp'$ matches the concatenated string xx' . Similarly, if $exp = (00|11)^*$ then $\Phi_{exp}(x) = 1$ if and only if x is of even length and $x_{2i} = x_{2i+1}$ for every $i < |x|/2$. We now turn to the formal definition of Φ_{exp} .

P The formal definition of Φ_{exp} is one of those definitions that is more cumbersome to write than to grasp. Thus it might be easier for you to first work it out on your own and then check that your definition matches what is written below.

Definition 9.2 — Matching a regular expression. Let exp be a regular expression. Then the function Φ_{exp} is defined as follows:

1. If $exp = \sigma$ then $\Phi_{exp}(x) = 1$ iff $x = \sigma$.
2. If $exp = (exp' | exp'')$ then $\Phi_{exp}(x) = \Phi_{exp'}(x) \vee \Phi_{exp''}(x)$ where \vee is the OR operator.
3. If $exp = (exp')(exp'')$ then $\Phi_{exp}(x) = 1$ iff there is some $x', x'' \in \Sigma^*$ such that x is the concatenation of x' and x'' and $\Phi_{exp'}(x') = \Phi_{exp''}(x'') = 1$.
4. If $exp = (exp')^*$ then $\Phi_{exp}(x) = 1$ iff there are is $k \in \mathbb{N}$ and some $x_0, \dots, x_{k-1} \in \Sigma^*$ such that x is the concatenation $x_0 \dots x_{k-1}$ and $\Phi_{exp'}(x_i) = 1$ for every $i \in [k]$.
5. Finally, for the edge cases Φ_{\emptyset} is the constant zero function, and Φ_{\cdot} is the function that only outputs 1 on the constant string.

We say that a function $F : \Sigma^* \rightarrow \{0, 1\}$ is *regular* if $F = \Phi_{exp}$ for some regular expression exp .⁴

For example let $\Sigma = \{a, b, c, d, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ and $F : \Sigma^* \rightarrow$

³ These are the regular expressions corresponding to accepting no strings, and accepting only the empty string respectively.

⁴ We use function notation in this book, but many other texts common talk about *languages*, which are sets of string. We say that a set $L \subseteq \Sigma^*$ (also known as a *language*) is *regular* if the corresponding function $F_L : \Sigma^* \rightarrow \{0, 1\}$ s.t. $F_L(x) = 1$ iff $x \in L$ is regular.

$\{0, 1\}$ be the function such that $F(x)$ outputs 1 iff x consists of one or more of the letters a - b followed by a sequence of one or more digits (without a leading zero). As shown by Eq. (9.2), the function $F : \Sigma^* \rightarrow \{0, 1\}$ is computable by a regular expression. For example the string $abc12078$ matches the expression of Eq. (9.2) since $\Phi_{a|b|c|d}(a) = 1$, $\Phi_{(a|b|c|d)^*}(bcd) = 1$, $\Phi_{(1|2|\dots|9)}(1) = 1$, and $\Phi_{(0|\dots|9)^*}(2078) = 1$.

P The definitions above are not inherently difficult, but are a bit cumbersome. So you should pause here and go over it again until you understand why it corresponds to our intuitive notion of regular expressions. This is important not just for understanding regular expressions themselves (which are used time and again in a great many applications) but also for getting better at understanding recursive definitions in general.

We can think of regular expressions as a type of “programming language”. That is, we can think of a regular expression exp over the alphabet Σ as a program that computes some function $\Phi : \Sigma^* \rightarrow \{0, 1\}$.⁵ But it turns out that the “halting problem” for these programs is easy: they always halt.

Theorem 9.3 — Regular expression always halt. For every set Σ and $exp \in (\Sigma \cup \{(\,), |, *, \emptyset, ""\})^*$, if exp is a valid regular expression over Σ then Φ_{exp} is a total function from Σ^* to $\{0, 1\}$. Moreover, there is an always halting NAND++ program P_{exp} that computes Φ_{exp} .

Proof Idea: The main idea behind the proof is to see that Definition 9.2 actually specifies a recursive algorithm for computing Φ_{exp} . The details are specified below. ★

Proof of Theorem 9.3. Definition 9.2 gives a way of recursively computing Φ_{exp} . The key observation is that in our recursive definition of regular expressions, whenever exp is made up of one or two expressions exp' , exp'' then these two regular expressions are *smaller* than exp , and eventually (when they have size 1) then they must correspond to the non-recursive case of a single alphabet symbol.

Therefore, we can prove the theorem by induction over the length m of exp (i.e., the number of symbols in the string exp , also denoted as $|exp|$). For $m = 1$, exp is a single alphabet symbol and the function Φ_{exp} is trivial. In the general case, for $m = |exp|$ we assume by the induction hypothesis that we have proven the theorem for $|exp| = 1, \dots, m - 1$. Then by the definition of regular expressions, exp is made up of one or two sub-expressions exp' , exp'' of length smaller

⁵ Regular expressions (and context free grammars, which we'll see below) are often thought of as *generative models* rather than computational ones, since their definition does not immediately give rise to a way to *decide* matches but rather to a way to generate matching strings by repeatedly choosing which rules to apply.

than m , and hence by the induction hypothesis we assume that $\Phi_{exp'}$ and $\Phi_{exp''}$ are total computable functions. But then we can follow the definition for the cases of concatenation, union, or the star operator to compute Φ_{exp} using $\Phi_{exp'}$ and $\Phi_{exp''}$. ■

9.2.1 Efficient matching of regular expressions (advanced, optional)

The proof of [Theorem 9.3](#) gives a recursive algorithm to evaluate whether a given string matches or not a regular expression. However, it turns out that there is a much more efficient algorithm to match regular expressions. One way to obtain such an algorithm is to replace this recursive algorithm with **dynamic programming**, using the technique of **memoization**.⁶ It turns out that the resulting dynamic program only requires maintaining a finite (independent of the input length) amount of state, and makes a single pass over its input. This means this algorithm is a **deterministic finite automaton (DFA)**. The relation of regular expressions with finite automata is a beautiful topic, on which we only touch upon in this text. See books such as [Sipser's](#), [Hopcroft, Motwani and Ullman](#), and [Kozen's](#).

We now prove the algorithmic result that regular expression matching can be done by a linear (i.e., $O(n)$) time algorithm, and moreover one that uses a constant (i.e., $O(1)$) amount of memory, and makes a single pass over its input. Since we have not yet covered the topics of time and space complexity, the reader might want to skip ahead at this point, and return to this theorem later.

Theorem 9.4 — DFA and regular expression equivalence. Let exp be a regular expression. Then there is an $O(n)$ time $O(1)$ space single pass algorithm (i.e., deterministic finite automaton) that computes Φ_{exp} .

More formally, there is an enhanced NAND++ program P that computes Φ_{exp} and moreover, P uses no array variable apart from **X**, **Xvalid**, **Y** and **Yvalid**, uses only the **i++** (`f○○`) operation (hence never goes back, only forward), and halts when **i** reaches beyond the length of the input.

Proof Idea: The idea is to first obtain a more efficient recursive algorithm for computing Φ_{exp} and then turing this recursive algorithm into a constant-space single-pass algorithm using the technique of *memoization*. In this technique we record in a table the results of every call to a function, and then if we make future calls with the same input, we retrieve the result from the table instead of re-computing it. This simple optimization can sometimes result in huge savings in running time.

For this case, we can define a recursive algorithm to compute Φ_{exp}

⁶ If you haven't taken yet an algorithms course, you might not know these techniques. This is OK; while the more efficient algorithm is crucial for the many practical applications of regular expressions, it is not of great importance to this course.

as follows: given $x \in \{0, 1\}^n$, if $n = 0$ then we output 1 if exp contains "", otherwise we output $\Phi_{exp[x_{n-1}]}(x_0 \cdots x_{n-1})$ where $exp[\sigma]$ is the result of "restricting" the output of exp to strings that end with σ . It can be shown that for every regular expression exp , we can construct such a regular expression $exp[\sigma]$ that would match a string x if and only if exp matches $x\sigma$. The resulting recursive algorithm runs in $O(n)$ time, we can then use memoization to make it into a single-pass constant space algorithm. Specifically, we will store a table of the (constantly many) expressions of length at most $|exp|$ that we need to deal with in the course of this algorithm, and iteratively for $i = 0, 1, \dots, n - 1$, compute whether or not each one of those expressions matches $x_0 \cdots x_{i-1}$. ★

Proof of Theorem 9.4. For a regular expression exp over an alphabet Σ and symbol $\sigma \in \Sigma$, we will define $exp[\sigma]$ to be a regular expression such that $exp[\sigma]$ matches a string x if and only if exp matches the string $x\sigma$. We can define $exp[\sigma]$ recursively as follows. For simplicity it will be a little more convenient to work with regular expressions where no sub-expression matches the empty string, except the explicit "" expression. That means that if we have an expression of the form $exp'exp''$ where exp'' can match the empty string, then we replace exp'' with an expression exp''' that matches all the strings that exp'' does except the empty string, and re-write the expression as $exp'|exp'exp'''$. Similarly, we will assume that all calls to the * operator match at least one occurrence, by enforcing that they are always of the form $exp|(exp*)exp$. (Once again, we can add an explicit expression for the empty string if we need to match it.) We call an expression that has the form above a *normal form expression*, and leave as an exercise to the reader to show that every expression exp has an equivalent expression that is in normal form.

Given a normal form expression exp , we transform it to $exp[\sigma]$ as follows:

1. If $exp = \tau$ for $\tau \in \Sigma$ then $exp[\sigma] = ""$ if $\tau = \sigma$ and $exp[\sigma] = \emptyset$ otherwise.
2. If $exp = exp'|exp''$ then $exp[\sigma] = exp'[\sigma]|exp''[\sigma]$.
3. If $exp = exp' exp''$ then $exp = exp' exp''[\sigma]$.⁷
4. If $exp = exp|(exp*)exp$ then $exp[\sigma] = exp[\sigma]|(exp*)(exp[\sigma])$.
5. If $exp = ""$ or $exp = \emptyset$ then $exp[\sigma] = \emptyset$.

⁷ Note that we use here the assumption that, because our expression is in normal form, exp'' does not match the empty string.

We leave it as an exercise to prove the following two claims:

Claim 1: For every $x \in \{0, 1\}^*$, $\Phi_{exp}(x\sigma) = 1$ if and only if $\Phi_{exp[\sigma]}(x) = 1$

Claim 2: For every normal form exp , $|exp[\sigma]| \leq |exp|$ where we denote by $|exp'|$ the number of symbols in the expression exp' .

Both claims can be proved by induction by following the recursive definition. Note that for Claim 2, we treat σ as a single symbol, and also simplify expressions of the form $exp\sigma$ to exp .

We can now define a recursive algorithm for computing Φ_{exp} :

Algorithm *MATCH*(exp, x):

Inputs: exp is normal form regular expression, $x \in \Sigma^n$ for some $n \in \mathbb{N}$.

1. If $x = \epsilon$ then return 1 iff exp has the form $exp = \sigma|exp'$ for some exp' . (For a normal-form expression, this is the only way it matches the empty string.)
2. Otherwise, return $MATCH(exp[x_{n-1}], x_0 \cdots x_{n-1})$.

Algorithm *MATCH* is a recursive algorithm that on input an expression exp and a string $x \in \{0, 1\}^n$, does some constant time computation and then calls itself on input some expression exp' smaller than exp and a string x of length $n - 1$. Thus, if we define $T(\ell, n)$ to be the running time of the algorithm on expressions of length at most ℓ and inputs of length n , then we see that this satisfies the equation $T(\ell, n) \leq T(\ell, n - 1) + C(\ell)$ (where the $C(\ell)$ denotes the time it takes to scan over an ℓ length expression exp to transform it to the expression $exp[\sigma]$.) Hence we see that for every constant ℓ , the running time would be $O(n)$ where the hidden constant in the O notation can depend on ℓ .

Moreover, since a regular expression over alphabet Σ is simply a string over the alphabet $\Sigma \cup \{(\cdot), |, *, \sigma, \emptyset\}$, to give a crude upper bound, there are at most $(|\Sigma| + 10)^\ell$ expressions over this alphabet of length at most ℓ . Now, instead of computing *MATCH* recursively, we can compute it iteratively as follows:

Algorithm *MATCH'*(exp, x):

Inputs: exp is normal form regular expression, $x \in \Sigma^n$ for some $n \in \mathbb{N}$. Let $\ell = |exp|$.

Define variables $v_{exp'}$ for every exp' of length at most ℓ . Initially, $v_{exp'} = 1$ if and only if exp' matches the empty string.

- For $i = 0, \dots, n - 1$ do the following:
- For every exp' of length at most ℓ :
 - Let $temp_{exp'} = v_{exp'[x_i]}$
- every exp' of length at most ℓ :
 - Let $v_{exp'} = temp_{exp'}$

Output v_{exp} .

This algorithm maintains the invariant that at the end of step i , the variable $v_{exp'}$ is equal if and only if exp' matches the string $x_0 \dots x_{i-1}$. Note that it only maintains a constant number of variables, and that it proceeds in one linear scan over the input, and so this proves the theorem. ■

9.3 LIMITATIONS OF REGULAR EXPRESSIONS

The fact that functions computed by regular expressions always halt is of course one of the reasons why they are so useful. When you make a regular expression search, you are guaranteed that you will get a result. This is why operating systems, for example, restrict you for searching a file via regular expressions and don't allow searching by specifying an arbitrary function via a general-purpose programming language. But this always-halting property comes at a cost. Regular expressions cannot compute every function that is computable by NAND++ programs. In fact there are some very simple (and useful!) functions that they cannot compute, such as the following:

Theorem 9.5 — Matching parenthesis. Let $\Sigma = \{ \langle, \rangle \}$ and $MATCHPAREN : \Sigma^* \rightarrow \{0, 1\}$ be the function that given a string of parenthesis, outputs 1 if and only if every opening parenthesis is matched by a corresponding closed one. Then there is no regular expression over Σ that computes $MATCHPAREN$.

Theorem 9.5 is a consequence of the following result known as the *pumping lemma*:

Theorem 9.6 — Pumping Lemma. Let exp be a regular expression. Then there is some number n_0 such that for every $w \in \{0, 1\}^*$ with $|w| > n_0$ and $\Phi_{exp}(w) = 1$, it holds that we can write $w = xyz$ where

$|y| \geq 1$, $|xy| \leq n_0$ and such that $\Phi_{exp}(xy^kz) = 1$ for every $k \in \mathbb{N}$.

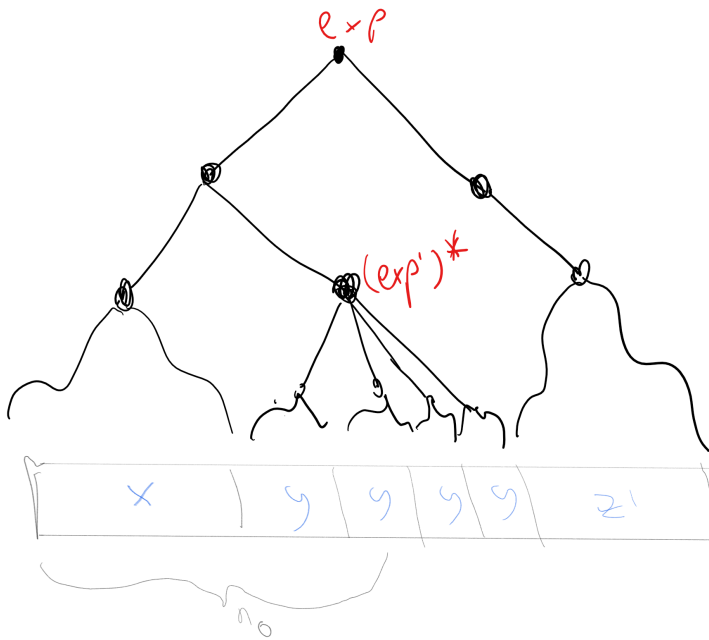


Figure 9.1: To prove the “pumping lemma” we look at a word w that is much larger than the regular expression exp that matches it. In such a case, part of w must be matched by some sub-expression of the form $(exp')^*$, since this is the only operator that allows matching words longer than the expression. If we look at the “leftmost” such sub-expression and define y^k to be the string that is matched by it, we obtain the partition needed for the pumping lemma.

Proof Idea: The idea behind the proof is very simple (see Fig. 9.1). If we let n_0 be, say, twice the number of symbols that are used in the expression exp , then the only way that there is some w with $|w| > n_0$ and $\Phi_{exp}(w) = 1$ is that exp contains the $*$ (i.e. star) operator and that there is a nonempty substring y of w that was matched by $(exp')^*$ for some sub-expression exp' of exp . We can now repeat y any number of times and still get a matching string. ★

P The pumping lemma is a bit cumbersome to state, but one way to remember it is that it simply says the following: “if a string matching a regular expression is long enough, one of its substrings must be matched using the $*$ operator”.

Proof of Theorem 9.6. To prove the lemma formally, we use induction on the length of the expression. Like all induction proofs, this is going to be somewhat lengthy, but at the end of the day it directly follows

the intuition above that *somewhere* we must have used the star operation. Reading this proof, and in particular understanding how the formal proof below corresponds to the intuitive idea above, is a very good way to get more comfort with inductive proofs of this form.

Our inductive hypothesis is that for an n length expression, $n_0 = 2n$ satisfies the conditions of the lemma. The base case is when the expression is a single symbol or that it is \emptyset or ϵ in which case the condition is satisfied just because there is no matching string of length more than one. Otherwise, exp is of the form **(a)** $exp'exp''$, **(b)**, $(exp')(exp'')$, **(c)** or $(exp'')^*$ where in all these cases the subexpressions have fewer symbols than exp and hence satisfy the induction hypothesis.

In case **(a)**, every string w matching exp must match either exp' or exp'' . In the former case, since exp' satisfies the induction hypothesis, if $|w| > n_0$ then we can write $w = xyz$ such that xy^kz matches exp' for every k , and hence this is matched by exp as well.

In case **(b)**, if w matches $(exp')(exp'')$. then we can write $w = w'w''$ where w' matches exp' and w'' matches exp'' . Again we split to subcases. If $|w'| > 2|exp'|$, then by the induction hypothesis we can write $w' = xyz$ of the form above such that xy^kz matches exp' for every k and then xy^kzw'' matches $(exp')(exp'')$. This completes the proof since $|xy| \leq 2|exp'|$ and so in particular $|xy| \leq 2(|exp'| + |exp''|) \leq 2|exp|$, and hence zw'' can be play the role of z in the proof. Otherwise, if $|w'| \leq 2|exp'|$ then since $|w|$ is larger than $2|exp|$ and $w = w'w''$ and $exp = exp'exp''$, we get that $|w'| + |w''| > 2(|exp'| + |exp''|)$. Thus, if $|w'| \leq 2|exp'|$ it must be that $|w''| > 2|exp''|$ and hence by the induction hypothesis we can write $w'' = xyz$ such that xy^kz matches exp'' for every k and $|xy| \leq 2|exp''|$. Therefore we get that $w'xy^kz$ matches $(exp')(exp'')$ for every k and since $|w'| \leq 2|exp'|$, $|w'xy| \leq 2(|exp'| + |exp''|)$ and this completes the proof since $w'x$ can play the role of x in the statement.

Now in the case **(c)**, if w matches $(exp'')^*$ then $w = w_0 \cdots w_t$ where w_i is a nonempty string that matches exp'' for every i . If $|w_0| > 2|exp''|$ then we can use the same approach as in the concatenation case above. Otherwise, we simply note that if x is the empty string, $y = w_0$, and $z = w_1 \cdots w_t$ then xy^kz will match $(exp'')^*$ for every k . ■

R Recursive definitions and inductive proofs When an object is *recursively defined* (as in the case of regular expressions) then it is natural to prove properties of such objects by *induction*. That is, if we want to prove that all objects of this type have property P , then it is natural to use an inductive steps that says that if o', o'', o''' etc have property P then so is an object o that is obtained by composing them.

Given the pumping lemma, we can easily prove [Theorem 9.5](#):

Proof of Theorem 9.5. Suppose, towards the sake of contradiction, that there is an expression exp such that $\Phi_{exp} = MATCHPAREN$. Let n_0 be the number from [Theorem 9.5](#) and let $w = \langle n_0 \rangle^{n_0}$ (i.e., n_0 left parenthesis followed by n_0 right parenthesis). Then we see that if we write $w = xyz$ as in [Theorem 9.5](#), the condition $|xy| \leq n_0$ implies that y consists solely of left parenthesis. Hence the string xy^2z will contain more left parenthesis than right parenthesis. Hence $MATCHPAREN(xy^2z) = 0$ but by the pumping lemma $\Phi_{exp}(xy^2z) = 1$, contradicting our assumption that $\Phi_{exp} = MATCHPAREN$. ■

The pumping lemma is a very useful tool to show that certain functions are *not* computable by a regular language. However, it is *not* an “if and only if” condition for regularity. There are non regular functions which still satisfy the conditions of the pumping lemma. To understand the pumping lemma, it is important to follow the order of quantifiers in [Theorem 9.6](#). In particular, the number n_0 in the statement of [Theorem 9.6](#) depends on the regular expression (in particular we can choose n_0 to be twice the number of symbols in the expression). So, if we want to use the pumping lemma to rule out the existence of a regular expression exp computing some function F , we need to be able to choose an appropriate w that can be arbitrarily large and satisfies $F(w) = 1$. This makes sense if you think about the intuition behind the pumping lemma: we need w to be large enough as to force the use of the star operator.

Solved Exercise 9.1 — Palindromes is not regular. Prove that the following function over the alphabet $\{0, 1, ;\}$ is not regular: $PAL(w) = 1$ if and only if $w = u; u^R$ where $u \in \{0, 1\}^*$ and u^R denotes u “reversed”: the string $u_{|u|-1} \dots u_0$.⁸ ■

Solution: We use the pumping lemma. Suppose towards the sake of contradiction that there is a regular expression exp computing PAL , and let n_0 be the number obtained by the pumping lemma ([Theorem 9.6](#)). Consider the string $w = 0^{n_0}; 0^{n_0}$. Since the reverse of the all zero string is the all zero string, $PAL(w) = 1$. Now, by the pumping lemma, if PAL is computed by exp , then we can write $w = xyz$ such that $|xy| \leq n_0$, $|y| \geq 1$ and $PAL(xy^kz) = 1$ for every $k \in \mathbb{N}$. In particular, it must hold that $PAL(xz) = 1$, but this is a contradiction, since $xz = 1^{n_0-|y|}; 0^{n_0}$ and so its two parts are not of the same length and in particular are not the reverse of one another. ■

⁸ The *Palindrome* function is most often defined without an explicit separator character $;$, but the version with such a separator is a bit cleaner and so we use it here. This does not make much difference, as one can easily encode the separator as a special binary string instead.

9.4 OTHER SEMANTIC PROPERTIES OF REGULAR EXPRESSIONS

Regular expressions are widely used beyond just searching. First, they are typically used to define *tokens* in various formalisms such as programming data description languages. But they are also used beyond it. One nice example is the recent work on the **NetKAT network programming language**. In recent years, the world of networking moved from fixed topologies to “software defined networks”, that are run by programmable switches that can implement policies such as “if packet is SSL then forward it to A, otherwise forward it to B”. By its nature, one would want to use a formalism for such policies that is guaranteed to always halt (and quickly!) and that where it is possible to answer semantic questions such as “does C see the packets moved from A to B” etc. The NetKAT language uses a variant of regular expressions to achieve that.

Such applications use the fact that, due to their restrictions, we can solve not just the halting problem for them, but also answer several other semantic questions as well, all of whom would not be solvable for Turing complete models due to Rice’s Theorem ([Theorem 8.7](#)). For example, we can tell whether two regular expressions are *equivalent*, as well as whether a regular expression computes the constant zero function.

Theorem 9.7 — Emptiness of regular languages is computable.. There is an algorithm that given a regular expression exp , outputs 1 if and only if Φ_{exp} is the constant zero function.

Proof Idea: The idea is that we can directly observe this from the structure of the expression. The only way it will output the constant zero function is if it has the form \emptyset or is obtained by concatenating \emptyset with other expressions. ★

Proof of Theorem 9.7. Define a regular expression to be “empty” if it computes the constant zero function. The algorithm simply follows the following rules:

- If an expression has the form σ or ϵ then it is not empty.
- If exp is not empty then $exp|exp'$ is not empty for every exp' .
- If exp is not empty then exp^* is not empty.
- If exp and exp' are both not empty then $exp exp'$ is not empty.
- \emptyset is empty.

- σ and ϵ are not empty.

Using these rules it is straightforward to come up with a recursive algorithm to determine emptiness. We leave verifying the details to the reader. ■

Theorem 9.8 — Equivalence of regular expressions is computable..

There is an efficient algorithm that on input two regular expressions exp, exp' , outputs 1 if and only if $\Phi_{exp} = \Phi_{exp'}$.

Proof Idea: [Theorem 9.7](#) above is actually a special case of [Theorem 9.8](#), since emptiness is the same as checking equivalence with the expression \emptyset . However we prove [Theorem 9.8](#) from [Theorem 9.7](#). The idea is that given exp and exp' , we will compute an expression exp'' such that $\Phi_{exp''}(x) = (\Phi_{exp}(x) \wedge \overline{\Phi_{exp'}(x)}) \vee (\overline{\Phi_{exp}(x)} \wedge \Phi_{exp'}(x))$ (where \bar{y} denotes the negation of y , i.e., $\bar{y} = 1 - y$). One can see that exp is equivalent to exp' if and only if exp'' is empty. To construct this expression, we need to show how given expressions exp and exp' , we can construct expressions $exp \wedge exp'$ and \overline{exp} that compute the functions $\Phi_{exp} \wedge \Phi_{exp'}$ and $\overline{\Phi_{exp}}$ respectively. (Computing the expression for $exp \vee exp'$ is straightforward using the $|$ operation of regular expressions.) Using De-Morgan's laws, it is enough to compute the negation, since the AND function can be expressed using OR and NOT. This is cumbersome but ultimately can be done, see details below. ★

Proof of [Theorem 9.8](#). The heart of the proof is the following claim:

Claim: For every regular expression exp over the alphabet Σ there is an expression \overline{exp} such that $\Phi_{\overline{exp}}(x) = 1 - \Phi_{exp}(x)$ for every $x \in \Sigma^*$.

Proof: We'll prove the claim for the case of $\Sigma = \{0, 1\}^*$. Extending this to other finite alphabets is straightforward. We prove the claim by recursion. We will use the ideas from the proof of [Theorem 9.4](#), and specifically we will assume that exp is in normal form, and also use the notation $exp[\sigma]$, as defined in that proof.

If $exp = \emptyset$ then $\overline{exp} = (0|1)^*$. Otherwise, we define $\overline{exp} = \overline{exp[0]0|exp[1]1}$ where $exp[\sigma]$ is the expression that matches x if and only if exp matches $x\sigma$. This gives us a recursive definition for \overline{exp} .⁹

Let $x \in \{0, 1\}^n$. We will prove by induction on n that exp matches x if and only if \overline{exp} does *not* match x . This is trivially true for the case that $n = 0$. Now if exp matches $x \in \{0, 1\}^n$, then letting $x' = x_0 \cdots x_{n-2}$ and $\sigma = x_{n-1}$, exp matches x if and only if $exp[\sigma]$

matches x' which happens if and only if $\overline{exp[\sigma]}$ does *not* match x' . But \overline{exp} matches $x'\sigma$ if and only if $\overline{exp[\sigma]}$ matches x' , hence completing the proof of the claim.

Once we have the claim, we can reduce checking equivalence to checking emptiness. For every two expressions exp and exp' we can define $exp \vee exp'$ to be simply the expression $exp|exp'$ and $exp \wedge exp'$ as $\overline{exp} \vee \overline{exp'}$. Now we can define

$$exp'' = (exp \wedge \overline{exp'}) \vee (\overline{exp} \wedge exp') \quad (9.3)$$

and verify that $\Phi_{exp''}$ is the constant zero function if and only if $\Phi_{exp}(x) = \Phi_{exp'}(x)$ for every $x \in \Sigma^*$. ■

9.5 CONTEXT FREE GRAMMARS

If you have ever written a program, you've experienced a *syntax error*. You might also have had the experience of your program entering into an *infinite loop*. What is less likely is that the compiler or interpreter entered an infinite loop when trying to figure out if your program has a syntax error.

When a person designs a programming language, they need to come up with a function $VALID : \{0, 1\}^* \rightarrow \{0, 1\}$ that determines the strings that correspond to valid programs in this language. The compiler or interpreter computes $VALID$ on the string corresponding to your source code to determine if there is a syntax error. To ensure that the compiler will always halt in this computation, language designers typically *don't* use a general Turing-complete mechanism to express the function $VALID$, but rather a restricted computational model. One of the most popular choices for such a model is *context free grammar*.

To explain context free grammars, let's begin with a canonical example. Let us try to define a function $ARITH : \Sigma^* \rightarrow \{0, 1\}$ that takes as input a string x over the alphabet $\Sigma = \{(\, , \, +, \, -, \, \times, \, \div, \, 0, \, 1, \, 2, \, 3, \, 4, \, 5, \, 6, \, 7, \, 8, \, 9\}$ and returns 1 if and only if the string x represents a valid arithmetic expression. Intuitively, we build expressions by applying an operation to smaller expressions, or enclosing them in parenthesis, where the "base case" corresponds to expressions that are simply numbers. A bit more precisely, we can make the following definitions:

- A *digit* is one of the symbols 0, 1, 2, 3, 4, 5, 6, 7, 8, 9.
- A *number* is a sequence of digits.¹⁰

⁹ To ensure this is well defined we need to use the fact that $exp[\sigma]$ is an expression of the same length or shorter than exp , and moreover there are no "cycles" in the form that there is no string $\sigma = \sigma_0 \cdots \sigma_t$ such that $exp[\sigma]$, defined as $exp[\sigma_0][\sigma_1] \cdots [\sigma_t]$ is equal to exp . We can observe the latter by looking a string x such that $\Phi_{exp}(x) = 1$ and x ends with the minimal number of repetitions of the string $\sigma_0 \cdots \sigma_t$. This x will satisfy $\Phi_{exp[\sigma]}(x) = 0$.

¹⁰ For simplicity we drop the condition that the sequence does not have a leading zero, though it is not hard to encode it in a context-free grammar as well.

- An *operation* is one of $+$, $-$, \times , \div
- An *expression* has either the form “*number*” or the form “*subexpression1 operation subexpression2*” or “(*subexpression*)”.

A context free grammar (CFG) is a formal way of specifying such conditions. We can think of a CFG as a set of rules to *generate* valid expressions. In the example above, the rule $expression \Rightarrow expression \times expression$ tells us that if we have built two valid expressions $exp1$ and $exp2$, then the expression $exp1 \times exp2$ is valid above.

We can divide our rules to “base rules” and “recursive rules”. The “base rules” are rules such as $number \Rightarrow 0$, $number \Rightarrow 1$, $number \Rightarrow 2$ and so on, that tell us that a single digit is a number. The “recursive rules” are rules such as $number \Rightarrow number digit$ that tell us that if we add a digit to a valid number then we still have a valid number. We now make the formal definition of context-free grammars:

Definition 9.9 — Context Free Grammar. Let Σ be some finite set. A *context free grammar (CFG) over Σ* is a triple (V, R, s) where V is a set disjoint from Σ of *variables*, R is a set of *rules*, which are pairs (v, z) (which we will write as $v \Rightarrow z$) where $v \in V$ and $z \in (\Sigma \cup V)^*$, and $s \in V$ is the starting rule.

■ **Example 9.10 — Context free grammar for arithmetic expressions..**

The example above of well-formed arithmetic expressions can be captured formally by the following context free grammar:

- The alphabet Σ is $\{ (,), +, -, \times, \div, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$
- The variables are $V = \{ expression, number, digit, operation \}$.
- The rules correspond the set R containing the following pairs:
 - $operation \Rightarrow +, operation \Rightarrow -, operation \Rightarrow \times, operation \Rightarrow \div$
 - $digit \Rightarrow 0, \dots, digit \Rightarrow 9$
 - $number \Rightarrow digit$
 - $number \Rightarrow digit number$
 - $expression \Rightarrow number$
 - $expression \Rightarrow expression operation expression$
 - $expression \Rightarrow (expression)$

- The starting variable is *expression*

There are various notations to write context free grammars in the literature, with one of the most common being **Backus-Naur form** where we write a rule of the form $v \Rightarrow a$ (where v is a variable and a is a string) in the form $\langle v \rangle := a$. If we have several rules of the form $v \mapsto a, v \mapsto b$, and $v \mapsto c$ then we can combine them as $\langle v \rangle := a|b|c$ (and this similarly extends for the case of more rules). For example, the Backus-Naur description for the context free grammar above is the following (using ASCII equivalents for operations):

```
operation := +|-|*|/
digit     := 0|1|2|3|4|5|6|7|8|9
number    := digit|digit number
expression := number|expression operation
           ↪ expression|(expression)
```

Another example of a context free grammar is the “matching parenthesis” grammar, which can be represented in Backus-Naur as follows:

```
match := ""|match match|(match)
```

You can verify that a string over the alphabet $\{ (,) \}$ can be generated from this grammar (where `match` is the starting expression and `""` corresponds to the empty string) if and only if it consists of a matching set of parenthesis.

9.5.1 Context-free grammars as a computational model

We can think of a CFG over the alphabet Σ as defining a function that maps every string x in Σ^* to 1 or 0 depending on whether x can be generated by the rules of the grammars. We now make this definition formally.

Definition 9.11 — Deriving a string from a grammar. If $G = (V, R, s)$ is a context-free grammar over Σ , then for two strings $\alpha, \beta \in (\Sigma \cup V)^*$ we say that β can be derived in one step from α , denoted by $\alpha \Rightarrow_G \beta$, if we can obtain β from α by applying one of the rules of G . That is, we obtain β by replacing in α one occurrence of the variable v with the string z , where $v \Rightarrow z$ is a rule of G .

We say that β can be derived from α , denoted by $\alpha \Rightarrow_G^* \beta$, if it can be derived by some finite number k of steps. That is, if there are $\alpha_1, \dots, \alpha_{k-1} \in (\Sigma \cup V)^*$, so that $\alpha \Rightarrow_G \alpha_1 \Rightarrow_G \alpha_2 \Rightarrow_G \dots \Rightarrow_G \alpha_{k-1} \Rightarrow_G \beta$.

We define the function computed by (V, R, s) to be the map $\Phi_{V,R,s} : \Sigma^* \rightarrow \{0, 1\}$ such that $\Phi_{V,R,s}(x) = 1$ iff $s \Rightarrow_G^* x$.

We say that $F : \Sigma^* \rightarrow \{0, 1\}$ is *context free* if $F = \Phi_{V,R,s}$ for some CFG (V, R, s) .¹¹

A priori it might not be clear that the map $\Phi_{V,R,s}$ is computable, but it turns out that we can in fact compute it. That is, the “halting problem” for context free grammars is trivial:

Theorem 9.12 — Context-free grammars always halt. For every CFG (V, R, s) over Σ , the function $\Phi_{V,R,s} : \Sigma^* \rightarrow \{0, 1\}$ is computable.

Proof. We only sketch the proof. It turns out that we can convert every CFG to an equivalent version that has the so called *Chomsky normal form*, where all rules either have the form $u \rightarrow vw$ for variables u, v, w or the form $u \rightarrow \sigma$ for a variable u and symbol $\sigma \in \Sigma$, plus potentially the rule $s \rightarrow \epsilon$ where s is the starting variable. (The idea behind such a transformation is to simply add new variables as needed, and so for example we can translate a rule such as $v \rightarrow u\sigma w$ into the three rules $v \rightarrow ur, r \rightarrow tw$ and $t \rightarrow \sigma$.)

Using this form we get a natural recursive algorithm for computing whether $s \Rightarrow_G^* x$ for a given grammar G and string x . We simply try all possible guesses for the first rule $s \rightarrow uv$ that is used in such a derivation, and then all possible ways to partition x as a concatenation $x = x'x''$. If we guessed the rule and the partition correctly, then this reduces our task to checking whether $u \Rightarrow_G^* x'$ and $v \Rightarrow_G^* x''$, which (as it involves shorter strings) can be done recursively. The base cases are when x is empty or a single symbol, and can be easily handled. ■

R **Parse trees** While we present CFGs as merely *deciding* whether the syntax is correct or not, the algorithm to compute $\Phi_{V,R,s}$ actually gives more information than that. That is, on input a string x , if $\Phi_{V,R,s}(x) = 1$ then the algorithm yields the sequence of rules that one can apply from the starting vertex s to obtain the final string x . We can think of these rules as determining a connected directed acyclic graph (i.e., a *tree*) with s being a source (or *root*) vertex and the sinks (or *leaves*) corresponding to the substrings of x that are obtained by the rules that do not have a variable in their second element. This tree is known as the *parse tree* of x , and often yields very useful information about the structure of x . Often the first step in a compiler or interpreter for a programming language is a *parser* that transforms the source into the parse tree (often known in this context as the **abstract syntax tree**). There are also tools that can automatically convert a description of

¹¹ As in the case of [Definition 9.2](#) we can also use *language* rather than *function* notation and say that a language $L \subseteq \Sigma^*$ is *context free* if the function F such that $F(x) = 1$ iff $x \in L$ is context free.

a context-free grammars into a *parser* algorithm that computes the parse tree of a given string. (Indeed, the above recursive algorithm can be used to achieve this, but there are much more efficient versions, especially for grammars that have **particular forms**, and programming language designers often try to ensure their languages have these more efficient grammars.)

9.5.2 The power of context free grammars

While we can (and people do) talk about context free grammars over any alphabet Σ , in the following we will restrict ourselves to $\Sigma = \{0, 1\}$. This is of course not a big restriction, as any finite alphabet Σ can be encoded as strings of some finite size. It turns out that context free grammars can capture every regular expression:

Theorem 9.13 — Context free grammars and regular expressions. Let exp be a regular expression over $\{0, 1\}$, then there is a CFG (V, R, s) over $\{0, 1\}$ such that $\Phi_{V,R,s} = \Phi_{exp}$.

Proof. We will prove this by induction on the length of exp . If exp is an expression of one bit length, then $exp = 0$ or $exp = 1$, in which case we leave it to the reader to verify that there is a (trivial) CFG that computes it. Otherwise, we fall into one of the following case: **case 1:** $exp = exp'exp''$, **case 2:** $exp = exp'|exp''$ or **case 3:** $exp = (exp')^*$ where in all cases exp', exp'' are shorter regular expressions. By the induction hypothesis have grammars (V', R', s') and (V'', R'', s'') that compute $\Phi_{exp'}$ and $\Phi_{exp''}$ respectively. By renaming of variables, we can also assume without loss of generality that V' and V'' are disjoint.

In case 1, we can define the new grammar as follows: we add a new starting variable $s \notin V \cup V'$ and the rule $s \mapsto s's''$. In case 2, we can define the new grammar as follows: we add a new starting variable $s \notin V \cup V'$ and the rules $s \mapsto s'$ and $s \mapsto s''$. Case 3 will be the only one that uses *recursion*. As before we add a new starting variable $s \notin V \cup V'$, but now add the rules $s \mapsto \epsilon$ (i.e., the empty string) and also add for every rule of the form $(s', \alpha) \in R'$ the rule $s \mapsto s\alpha$ to R .

We leave it to the reader as (again a very good!) exercise to verify that in all three cases the grammars we produce capture the same function as the original expression. ■

It turns out that CFG's are strictly more powerful than regular expressions. In particular, as we've seen, the "matching parenthesis" function *MATCHPAREN* can be computed by a context free grammar, whereas, as shown in [Theorem 9.5](#), it cannot be computed by regular expressions. Here is another example:

Solved Exercise 9.2 — Context free grammar for palindromes. Let $PAL : \{0, 1, ;\}^* \rightarrow \{0, 1\}$ be the function defined in [Solved Exercise 9.1](#) where $PAL(w) = 1$ iff w has the form $u;u^R$. Then PAL can be computed by a context-free grammar ■

Solution: A simple grammar computing PAL can be described using Backus–Naur notation:

```
start      := ; | 0 start 0 | 1 start 1
```

One can prove by induction that this grammar generates exactly the strings w such that $PAL(w) = 1$. ■

A more interesting example is computing the strings of the form $u;v$ that are *not* palindromes:

Solved Exercise 9.3 — Non palindromes. Prove that there is a context free grammar that computes $NPAL : \{0, 1, ;\}^* \rightarrow \{0, 1\}$ where $NPAL(w) = 0$ if $w = u;v$ but $v \neq u^R$. ■

Solution: Using Backus–Naur notation we can describe such a grammar as follows

```
palindrome := ; | 0 palindrome 0 | 1
           ↪ palindrome 1
different  := 0 palindrome 1 | 1 palindrome
           ↪ 0
start      := different | 0 start | 1 start
           ↪ | start 0 | start 1
```

In words, this means that we can characterize a string w such that $NPAL(w) = 1$ as having the following form

$$w = \alpha b u ; u^R b' \beta \quad (9.4)$$

where α, β, u are arbitrary strings and $b \neq b'$. Hence we can generate such a string by first generating a palindrome $u;u^R$ (palindrome variable), then adding either 0 on the right and 1 on the left to get something that is *not* a palindrome (different variable), and then we can add arbitrary number of 0's and 1's on either end (the start variable). ■

9.5.3 Limitations of context-free grammars (optional)

Even though context-free grammars are more powerful than regular expressions, there are some simple languages that are *not* captured by context free grammars. One tool to show this is the context-free grammar analog of the “pumping lemma” ([Theorem 9.6](#)):

Theorem 9.14 — Context-free pumping lemma. Let (V, R, s) be a CFG over Σ , then there is some $n_0 \in \mathbb{N}$ such that for every $x \in \Sigma^*$ with $|\sigma| > n_0$, if $\Phi_{V,R,s}(x) = 1$ then $x = abcde$ such that $|b| + |c| + |d| \leq n_1$, $|b| + |d| \geq 1$, and $\Phi_{V,R,s}(ab^kcd^ke) = 1$ for every $k \in \mathbb{N}$.

P The context-free pumping lemma is even more cumbersome to state than its regular analog, but you can remember it as saying the following: “If a long enough string is matched by a grammar, there must be a variable that is repeated in the derivation.”

Proof of Theorem 9.14. We only sketch the proof. The idea is that if the total number of symbols in the rules R is k_0 , then the only way to get $|x| > k_0$ with $\Phi_{V,R,s}(x) = 1$ is to use *recursion*. That is there must be some $v \in V$ such that by a sequence of rules we are able to derive from v the value bvd for some strings $b, d \in \Sigma^*$ and then further on derive from v the string $c \in \Sigma^*$ such that bcd is a substring of x . If try to take the minimal such v then we can ensure that $|bcd|$ is at most some constant depending on k_0 and we can set n_0 to be that constant ($n_0 = 10|R|k_0$ will do, since we will not need more than $|R|$ applications of rules, and each such application can grow the string by at most k_0 symbols). Thus by the definition of the grammar, we can repeat the derivation to replace the substring bcd in x with b^kcd^k for every $k \in \mathbb{N}$ while retaining the property that the output of $\Phi_{V,R,s}$ is still one. ■

Using [Theorem 9.14](#) one can show that even the simple function $F(x) = 1$ iff $x = ww$ for some $w \in \{0, 1\}^*$ is not context free. (In contrast, the function $F(x) = 1$ iff $x = ww^R$ for $w \in \{0, 1\}^*$ where for $w \in \{0, 1\}^n$, $w^R = w_{n-1}w_{n-2} \cdots w_0$ is context free, can you see why?.)

Solved Exercise 9.4 — Equality is not context-free. Let $EQ : \{0, 1, ;\}^* \rightarrow \{0, 1\}$ be the function such that $F(x) = 1$ if and only if $x = u;u$ for some $u \in \{0, 1\}^*$. Then EQ is not context free. ■

Solution: We use the context-free pumping lemma. Suppose towards the sake of contradiction that there is a grammar G that computes EQ , and let n_0 be the constant obtained from [Theorem 9.14](#). Consider the string $x = 1^{n_0}0^{n_0};1^{n_0}0^{n_0}$, and write it as $x = abcde$ as per [Theorem 9.14](#), with $|bcd| \leq n_0$ and with $|b| + |d| \geq 1$. By [Theorem 9.14](#), it should hold that $EQ(ace) = 1$. However, by case analysis this can be shown to be a contradiction.

First of all, unless b is on the left side of the ; separator and d is on the right side, dropping b and d will definitely make the two parts different. But if it is the case that b is on the left side and d is on the right side, then by the condition that $|bcd| \leq n_0$ we know that b is a string of only zeros and d is a string of only ones. If we drop b and d then since one of them is non empty, we get that there are either less zeroes on the left side than on the right side, or there are less ones on the right side than on the left side. In either case, we get that $EQ(ace) = 0$, obtaining the desired contradiction. ■

9.6 SEMANTIC PROPERTIES OF CONTEXT FREE LANGUAGES

As in the case of regular expressions, the limitations of context free grammars do provide some advantages. For example, emptiness of context free grammars is decidable:

Theorem 9.15 — Emptiness for CFG's is decidable. There is an algorithm that on input a context-free grammar G , outputs 1 if and only if Φ_G is the constant zero function.

Proof Idea: The proof is easier to see if we transform the grammar to Chomsky Normal Form as in [Theorem 9.12](#). Given a grammar G , we can recursively define a non-terminal variable v to be *non empty* if there is either a rule of the form $v \Rightarrow \sigma$, or there is a rule of the form $v \Rightarrow uw$ where both u and w are non empty. Then the grammar is non empty if and only if the starting variable s is non-empty. ★

Proof of Theorem 9.15. We assume that the grammar G in Chomsky Normal Form as in [Theorem 9.12](#). We consider the following procedure for marking variables as “non empty”:

1. We start by marking all variables v that are involved in a rule of the form $v \Rightarrow \sigma$ as non empty.
2. We then continue to mark v as non empty if it is involved in a rule of the form $v \Rightarrow uw$ where u, w have been marked before.

We continue this way until we cannot mark any more variables. We then declare that the grammar is empty if and only if s has not been marked. To see why this is a valid algorithm, note that if a variable v has been marked as “non empty” then there is some string $\alpha \in \Sigma^*$ that can be derived from v . On the other hand, if v has not been marked, then every sequence of derivations from v will always have a variable that has not been replaced by alphabet symbols. Hence in particular Φ_G is the all zero function if and only if the starting variable s is not marked “non empty”. ■

9.6.1 Uncomputability of context-free grammar equivalence (optional)

By analogy to regular expressions, one might have hoped to get an algorithm for deciding whether two given context free grammars are equivalent. Alas, no such luck. It turns out that the equivalence problem for context free grammars is *uncomputable*. This is a direct corollary of the following theorem:

Theorem 9.16 — Fullness of CFG's is uncomputable.. For every set Σ , let $CFGFULL_{\Sigma}$ be the function that on input a context-free grammar G over Σ , outputs 1 if and only if G computes the constant 1 function. Then there is some finite Σ such that $CFGFULL_{\Sigma}$ is uncomputable.

Theorem 9.16 immediately implies that equivalence for context-free grammars is uncomputable, since computing “fullness” of a grammar G over some alphabet $\Sigma = \{\sigma_0, \dots, \sigma_{k-1}\}$ corresponds to checking whether G is equivalent to the grammar $s \Rightarrow \#|s\sigma_0|\dots|s\sigma_{k-1}$. Note that **Theorem 9.16** and **Theorem 9.15** together imply that context-free grammars, unlike regular expressions, are *not* closed under complement. (Can you see why?)

Proof Idea: We prove the theorem by reducing from the Halting problem. To do that we use the notion of *configurations* of NAND++ programs, as defined in **Definition 7.12**. Recall that a *configuration* of a program P is a string s over some alphabet Σ that encodes all the information about the program in current iteration

We will let $\#$ and $;$ be some “separator characters” and then we can encode the *computation history* of a NAND++ program P on some input x by a sequence $s_0\#s_1^R; s_1\#s_2^R; \dots; s_{t-3}\#s_{t-2}^R; s_{t-2}\#s_{t-1}^R$ where s_0 is the starting configuration (on the input x) and s_{t-1} is a final configuration. We will show that for every NAND++ program P there is some finite alphabet Σ (containing $\#$ and $;$) and a grammar G_P that generates a string $\tau \in \Sigma^*$ if and only if τ does not encode a valid computation history of P on the input 0. Hence P halts on the empty input if and only if there is some τ that G_M can not generate, thus proving the theorem. ★

Proof of Theorem 9.16. We only sketch the proof. We will show that if we can compute $CFGFULL$ then we can solve *HALT*, which has been proven uncomputable in **Theorem 8.3**. Let P be an input program for *HALT* and x an input for P . We assume without loss of generality that P is well formed with a array variables and b scalar variables as in **Definition 7.12**. To prove the theorem we will show a grammar that can generate all strings over the alphabet

$\Sigma = \{0, 1\}^a \cup \{0, 1\}^{a+b} \cup \{\#, ;\}$ *except* those that correspond to valid histories of the form

$$s_0 \# s_1^R; s_1 \# s_2^R; \dots; s_{t-3} \# s_{t-2}^R; s_{t-2} \# s_{t-1}^R \quad (9.5)$$

where s_0 is a starting configuration, s_{t-1} is an ending configuration, and $s_{j+1} = \text{NEXT}_P(s_j)$ for every $j \in \{0, \dots, t-2\}$.

We will not spell out the full details but the main ideas are the following:

- Checking that a state is a valid starting or ending state can be done via a regular expression, and hence both this condition and its negation can be captured by context free grammars.
- We can also write a regular expression for the negation of the expression $((0|1)^* \# (0|1)^*)^*$ and so accept all strings that are not composed of blocks of this form.
- As we've seen in [Solved Exercise 9.3](#), we can find a context free grammar for all the strings *not* of the form $u; u^R$ (see .ref) and so can ensure our grammar accepts all string that contains such a “non matching” block.
- Now if a string is of the form $s_0 \# s_1^R; s_1 \# s_2^R; \dots; s_{t-3} \# s_{t-2}^R; s_{t-2} \# s_{t-1}^R$ then checking that it fails to be a valid computation history amounts to verifying that it contains a block of the form $s^R; s'$ such that s' is *not* equal to $\text{NEXT}_P(s)$. Since NEXT_P only modifies s in at most three locations, this can be done by extending the ideas in [Solved Exercise 9.3](#).

The above allows us to translate for every NAND++ program P , the question of whether P halts on zero to the question of whether a grammar G_P generates all strings. The only caveat is that the alphabet of our grammar grows with P , and so this does not necessarily show that CFGFULL_Σ is uncomputable for some *fixed* alphabet Σ . However, we can use the following observation: the Halting problem is uncomputable even if we restrict attention to the single universal NAND++ program U , as solving $\text{HALT}(P, x)$ is the same as solving $\text{HALT}(U, P \circ x)$ where \circ denotes concatenation and we identify P with its (prefix-free) encoding as a string. This shows that if U is a universal NAND++ program with a array variables and b scalar variables, then CFGFULL_Σ is uncomputable where $\Sigma = \{0, 1\}^a \cup \{0, 1\}^{a+b} \cup \{\#, ;\}$. ■

9.7 SUMMARY OF SEMANTIC PROPERTIES FOR REGULAR EXPRESSIONS AND CONTEXT-FREE GRAMMARS

To summarize, we can often trade *expressiveness* of the model for *amenability to analysis*. If we consider computational models that are

not Turing complete, then we are sometimes able to bypass Rice's Theorem and answer certain semantic questions about programs in such models. Here is a summary of some of what is known about semantic questions for the different models we have seen.

Model	Problem	Halting	Emptiness	Equivalence
Regular Expressions		Decidable	Decidable	Decidable
Context Free Grammars		Decidable	Decidable	Undecidable
Turing complete models		Undecidable	Undecidable	Undecidable

R **Unrestricted Grammars (optional)** The reason we call context free grammars "context free" is because if we have a rule of the form $v \mapsto a$ it means that we can always replace v with the string a , no matter the *context* in which v appears. More generally, we might want to consider cases where our replacement rules depend on the context.

This gives rise to the notion of *general grammars* that allow rules of the form $a \Rightarrow b$ where both a and b are strings over $(V \cup \Sigma)^*$. The idea is that if, for example, we wanted to enforce the condition that we only apply some rule such as $v \mapsto 0w1$ when v is surrounded by three zeroes on both sides, then we could do so by adding a rule of the form $000v000 \mapsto 0000w1000$ (and of course we can add much more general conditions). Alas, this generality comes at a cost - these general grammars are Turing complete and hence their halting problem is undecidable.

✓ **Lecture Recap**

- The uncomputability of the Halting problem for general models motivates the definition of restricted computational models.
- In some restricted models we can answer *semantic* questions such as: does a given program terminate, or do two programs compute the same function?
- *Regular expressions* are a restricted model of computation that is often useful to capture tasks of string matching. We can test efficiently whether an expression matches a string, as well as answer questions such as Halting and Equivalence.
- *Context free grammars* is a stronger, yet still not Turing complete, model of computation. The halting problem for context free grammars is computable, but equivalence is not computable.

9.8 EXERCISES

R **Disclaimer** Most of the exercises have been written in the summer of 2018 and haven't yet been fully debugged. While I would prefer people do not post online solutions to the exercises, I would greatly appreciate if you let me know of any bugs. You can do so by posting a [GitHub issue](#) about the exercise, and optionally complement this with an email to me with more details about the attempted solution.

9.9 BIBLIOGRAPHICAL NOTES

12

9.10 FURTHER EXPLORATIONS

Some topics related to this chapter that might be accessible to advanced students include: (to be completed)

9.11 ACKNOWLEDGEMENTS

¹² TODO: Add letter of Christopher Strachey to the editor of The Computer Journal. Explain right order of historical achievements. Talk about intuitionistic, logicist, and formalist approaches for the foundations of mathematics. Perhaps analogy to veganism. State the full Rice's Theorem and say that it follows from the same proof as in the exercise.