

Languages and Computation (COMP2012)  
Lecture notes  
Spring 2023

Thorsten Altenkirch, Venanzio Capretta, and Henrik Nilsson

January 28, 2023

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Example: Valid Java programs . . . . .	1
1.2	Example: The halting problem . . . . .	2
1.3	Example: The $\lambda$ -calculus . . . . .	4
1.4	P versus NP . . . . .	4
<b>2</b>	<b>Formal Languages</b>	<b>6</b>
2.1	Exercises . . . . .	8
<b>3</b>	<b>Finite Automata</b>	<b>10</b>
3.1	Deterministic finite automata . . . . .	10
3.1.1	What is a DFA? . . . . .	10
3.1.2	The language of a DFA . . . . .	12
3.2	Nondeterministic finite automata . . . . .	13
3.2.1	What is an NFA? . . . . .	13
3.2.2	The language accepted by an NFA . . . . .	15
3.2.3	The subset construction . . . . .	18
3.2.4	Correctness of the subset construction . . . . .	22
3.3	Exercises . . . . .	23
<b>4</b>	<b>Regular Expressions</b>	<b>26</b>
4.1	What are regular expressions? . . . . .	26
4.2	The meaning of regular expressions . . . . .	27
4.3	Algebraic laws . . . . .	30
4.4	Translating regular expressions into NFAs . . . . .	31
4.5	Summing up . . . . .	41
4.6	Exercises . . . . .	42
<b>5</b>	<b>Minimization of Finite Automata</b>	<b>43</b>
5.1	The table-filling algorithm . . . . .	43
5.2	Example of DFA minimization using the table-filling algorithm . . . . .	44

<b>6</b>	<b>Disproving Regularity</b>	<b>47</b>
6.1	The pumping lemma . . . . .	47
6.2	Applying the pumping lemma . . . . .	48
6.3	Exercises . . . . .	49
<b>7</b>	<b>Context-Free Grammars</b>	<b>50</b>
7.1	What are context-free grammars? . . . . .	50
7.2	The meaning of context-free grammars . . . . .	52
7.3	The relation between regular and context-free languages . . . . .	53
7.4	Derivation trees . . . . .	54
7.5	Ambiguity . . . . .	56
7.6	Applications of context-free grammars . . . . .	59
7.7	Exercises . . . . .	61
<b>8</b>	<b>Transformations of context-free grammars</b>	<b>63</b>
8.1	Equivalence of context-free grammars . . . . .	63
8.2	Elimination of useless productions . . . . .	63
8.3	Substitution . . . . .	64
8.4	Left factoring . . . . .	65
8.5	Disambiguating context-free grammars . . . . .	65
8.6	Elimination of left recursion . . . . .	67
8.7	Exercises . . . . .	71
<b>9</b>	<b>Pushdown Automata</b>	<b>72</b>
9.1	What is a pushdown automaton? . . . . .	72
9.2	How does a PDA work? . . . . .	73
9.3	The language of a PDA . . . . .	74
9.4	Deterministic PDAs . . . . .	75
9.5	Context-free grammars and push-down automata . . . . .	76
<b>10</b>	<b>Recursive-Descent Parsing</b>	<b>78</b>
10.1	What is parsing? . . . . .	78
10.2	Parsing strategies . . . . .	78
10.3	Basics of recursive-descent parsing . . . . .	79
10.4	Handling choice . . . . .	82
10.5	Recursive-descent parsing and left-recursion . . . . .	85
10.6	Predictive parsing . . . . .	85
10.6.1	First and follow sets . . . . .	87
10.6.2	LL(1) grammars . . . . .	87
10.6.3	Nullable nonterminals . . . . .	88
10.6.4	Computing first sets . . . . .	89
10.6.5	Computing follow sets . . . . .	90
10.6.6	Implementing a predictive parser . . . . .	91
10.6.7	LL(1), left-recursion, and ambiguity . . . . .	93
10.6.8	Satisfying the LL(1) conditions . . . . .	95
10.7	Beyond hand-written parsers: use parser generators . . . . .	96
10.8	Exercises . . . . .	97

<b>11 Turing Machines</b>	<b>99</b>
11.1 What is a Turing machine?	99
11.2 Grammars and context-sensitivity	102
11.3 The halting problem	103
11.4 Recursive and recursively enumerable sets	104
11.5 Back to Chomsky	107
11.6 Exercises	108
<b>12 <math>\lambda</math>-Calculus</b>	<b>110</b>
12.1 Syntax of $\lambda$ -calculus	110
12.2 Church numerals	113
12.3 Other data structures	114
12.4 Confluence	115
12.5 Recursion	116
12.6 The universality of $\lambda$ -calculus	117
12.7 Exercises	118
<b>13 Algorithmic Complexity</b>	<b>119</b>
13.1 The Satisfiability Problem	120
13.2 Time Complexity	120
13.3 $\mathcal{NP}$ -completeness	122
13.4 Exercises	123

## List of exercises

Exercise 2.1 t	8
Exercise 2.2 t	9
Exercise 2.3 t	9
Exercise 2.4 t	9
Exercise 3.1 t	23
Exercise 3.2 t	23
Exercise 3.3 t	23
Exercise 3.4 t	23
Exercise 3.5 t	24
Exercise 3.6 t	24
Exercise 3.7 t	25
Exercise 4.1 t	42
Exercise 4.2 t	42
Exercise 4.3 t	42
Exercise 4.4 t	42
Exercise 6.1 t	49
Exercise 7.1 t	61
Exercise 7.2 t	61
Exercise 7.3 t	61
Exercise 7.4 t	62
Exercise 8.1 t	71
Exercise 8.2 t	71
Exercise 8.3 t	71
Exercise 10.1 t	97

Exercise 10.2 t	97
Exercise 11.1 t	108
Exercise 11.2 t	108
Exercise 11.3 t	109
Exercise 12.1 t	118
Exercise 12.2 t	118
Exercise 13.1 t	123
Exercise 13.2 t	123
Exercise 13.3 t	124

## 1 Introduction

This module is about two fundamental notions in computer science, languages and computation, and how they are related. Specific topics include:

- Automata Theory
- Formal Languages
- Models of Computation
- Complexity Theory

The module starts with an investigation of classes of formal languages and related abstract machines, considers practical uses of this theory such as parsing, and finishes with a discussion on what computation is, what can and cannot be computed at all, and what can be computed efficiently, including famous results such as the Halting Problem and open problems such as P versus NP.

### 1.1 Example: Valid Java programs

Consider the following Java fragment:

```
class Foo {
    int n;
    void printNSqrd() {
        System.out.println(n * n);
    }
}
```

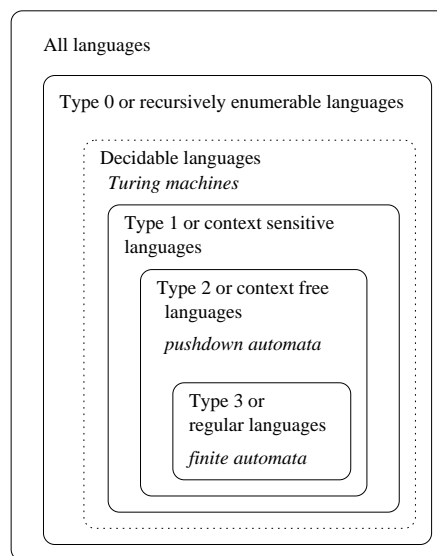
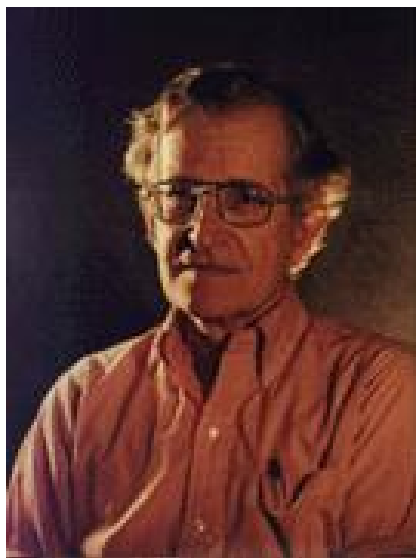
As written using a text editor or as stored in a file, it is just a string of characters. But not any string is a valid Java program. For example, Java uses specific keywords, have rules for what identifiers must look like, and requires proper nesting, such as a definition of a method inside a definition of a class.

This raises a number of questions:

- How to describe the set of strings that are valid Java programs?
- Given a string, how to determine if it is a valid Java program or not?
- How to recover the structure of a Java program from a “flat” string?

To answer such questions, we will study *regular expressions* and *grammars* to give precise descriptions of languages, and various kinds of *automata* to decide if a string belongs to a language or not. We will also consider how to systematically derive programs that efficiently answer this type of questions, drawing directly from the theory. Such programs are key parts of compilers, web browsers and web servers, and in fact of any program that uses structured text in one way or another.

A little bit of history. Context-free grammars were invented by American linguist, philosopher, and cognitive scientist *Noam Chomsky* (1928–) in an attempt to describe natural languages formally. He also introduced the *Chomsky Hierarchy* which classifies grammars and languages and their descriptive power:



## 1.2 Example: The halting problem

Consider the following program. Does it terminate for all values of  $n \geq 1$ ?

```
while (n > 1) {
    if even(n) {
        n = n / 2;
    } else {
        n = n * 3 + 1;
    }
}
```

This is not as easy to answer as it might first seem. Say we start with  $n = 7$ , for example:

7, 22, 11, 34, 17, 52, 26, 13, 40, 20, 10, 5, 16, 8, 4, 2, 1

Note how the numbers both increase and decrease in a way that is very hard to describe, which is exactly why it is so hard to analyse this program. The sequence involved is known as the *hailstone sequence*, and *Collatz conjecture* says that the number 1 will always be reached. And, in fact, for all numbers that have been

tried (all numbers up to  $2^{60}!$ ), the sequence does indeed terminate. But so far, no one has been able to *prove* that it always will! The famous mathematician Paul Erdős even said: “Mathematics may not be ready for such problems.” (See Collatz conjecture, Wikipedia.)

The following important decidability result should then perhaps not come as a total surprise:

It is impossible to write a program that decides if another, *arbitrary*, program terminates (halts) or not.

This is known as the *Halting Problem* and it is thus one example of an *undecidable* problem: the answer cannot be determined mechanically in general<sup>1</sup>.

The undecidability of the Halting Problem was first proved by British mathematician, logician, and computer scientist *Alan Turing* (1912–1954):



Turing proved this result using *Turing Machines*, an abstract model of computation that he introduced in 1936 to give a precise definition of what problems are “effectively calculable” (can be solved mechanically). Turing was further instrumental in the success of British code breaking efforts during World War II and is also famous for the Turing test to decide if a machine exhibits intelligent behaviour equivalent to, or indistinguishable from, that of a human. Andrew Hodges has written a very good biography of Turing: *Alan Turing: the Enigma* (<http://www.turing.org.uk/turing/>).

### 1.3 Example: The $\lambda$ -calculus

The  $\lambda$ -calculus is a theory of pure functions. It is very simple, having only two constructs: definition and application of functions. The following is an example of a  $\lambda$ -calculus term:

$$(\lambda x.x)(\lambda y.y)$$

---

<sup>1</sup>This does not mean that it is impossible to compute the answer for every instance of such a problem. On the contrary, in specific cases, the answer can often be computed very easily, and programs that attempt to solve undecidable problems can be very useful. But if we do write such a program, we must necessarily be prepared to give up one way or another with a “don’t know” answer.

Like Turing machines, the  $\lambda$ -calculus is a universal model of computation. It was introduced by American mathematician and logician *Alonzo Church* (1903–1995), also in 1936, a little earlier than Turing’s machine:



Alan Turing subsequently became a PhD student of Alonzo Church. They proved that Turing machines and the  $\lambda$ -calculus are equivalent in terms of computational expressiveness. In fact, all proposed universal models of computation to date have proved to be equivalent in that sense. This is captured by the *Church-Turing thesis*:

What is effectively calculable is exactly what can be computed by a Turing machine.

Functional programming languages, like Haskell, and many proof assistants implement (variations of) the  $\lambda$ -calculus. This is thus an example of a theory with very direct and practical applications.

## 1.4 P versus NP

Here is a seemingly innocuous question:

Can every problem whose solution can be *checked* quickly by a computer also be *solved* quickly by a computer?”

“Quickly” here means in time proportional to a *polynomial* in the size of the problem. Whether or not this is the case is known as the *P versus NP problem* and it is likely the most famous open problem in computer science, dating back to the 1950s. Here, “P” refers to the class of problems that can be solved in polynomial time, while “NP” refers to problems that can be solved in nondeterministic polynomial time, and the question is thus whether these two classes of problems actually are the same, or  $P = NP$ .

There is an abundance of important problems where solutions can be checked quickly, but where the best *known* algorithm for finding a solution is *exponential* in the size of the problem.

As an example, here is one, *the Subset Sum Problem*: Does some nonempty subset of given set of integers sum to zero? For example, given  $\{3, -2, 8, -5, 4, 9\}$ , the nonempty subset  $\{-5, -2, 3, 4\}$  sums to 0.

It is easy to check a proposed solution: just add all the numbers. If the initial set contains  $n$  integers, any proposed solution (being a subset) contains at most  $n$  integers, so we can sum all the elements with at most  $n$  additions meaning the total time taken is proportional to  $n$  (assuming addition is a constant time operation).

However, for *finding* a solution, no better way is known than essentially checking each possible subset one after another. As there are  $2^n$  subsets of a set with  $n$  elements, this means finding a solution takes exponential time.

Whether or not there is a better way to solve the Subset Sum Problem might not seem particularly important, but if it were the case that  $P = NP$ , then this would have monumental practical implications. For example, public key cryptography, on which pretty much all secure Internet communication, such as HTTPS, hinges, would no longer provide adequate security, and the entire Internet security infrastructure would have to be redesigned and reimplemented. The question here is if it is possible to quickly find the prime factors of (very) large numbers. As long as that is *not* the case, public key cryptography is considered secure.



## 2 Formal Languages

In this course we will use the terms *language* and *word* in a different way than in everyday language:

- A *language* is a set of words.
- A *word* is a sequence, or string, of symbols.

We will write  $\epsilon$  for the *empty* word; i.e., a sequence of length 0.

This leaves us with the question: what is a symbol? The answer is: anything, but it has to come from an alphabet  $\Sigma$  that is a finite set. A common (and important) instance is  $\Sigma = \{0, 1\}$ . Note that  $\epsilon$  will never be a symbol to avoid confusion.

Mathematically we say: Given an alphabet  $\Sigma$  we define the set  $\Sigma^*$  as set of words (or sequences) over  $\Sigma$ : the empty word  $\epsilon \in \Sigma^*$  and given a symbol  $x \in \Sigma$  and a word  $w \in \Sigma^*$  we can form a new word  $xw \in \Sigma^*$ . These are all the ways elements on  $\Sigma^*$  can be constructed (this is called an *inductive definition*). This unary  $*$ -operator is known as the *Kleene star* (or *Kleene operator* or *Kleene closure*).

With  $\Sigma = \{0, 1\}$ , typical elements of  $\Sigma^*$  are 0010, 00000000,  $\epsilon$ . Note, that we only write  $\epsilon$  if it appears on its own, instead of  $00\epsilon$  we just write 00.

Note further that  $\Sigma^*$  by definition is *always* nonempty as the empty word  $\epsilon$  belongs to  $\Sigma^*$  for *any* alphabet  $\Sigma$ , including  $\Sigma = \emptyset$ . Moreover, for any nonempty alphabet  $\Sigma$ ,  $\Sigma^*$  is an infinite set.

It is important to realise that although there are infinitely many words over a nonempty alphabet  $\Sigma$ , each word has a finite length. At first this may seem strange: how can it be that all elements of a set with infinitely many elements can be finite? A good way to think of an infinite set is as a *process* that can generate a new element whenever we need one, as many times as we like<sup>2</sup>. But each such element can obviously be of finite size as we at any point in time will only have asked for finitely many elements. Conversely, if we make a set containing a single (notionally) “infinite” element, such as a number  $\infty$  larger than any number except itself, or an infinitely long string, that does not make the set itself infinite: it would still contain exactly one element.

We can now define the notion of a language  $L$  over an alphabet  $\Sigma$  precisely:  $L \subseteq \Sigma^*$  or equivalently  $L \in \mathcal{P}(\Sigma^*)$ <sup>3</sup>.

Here are some informal examples of languages:

- The set  $\{0010, 00000000, \epsilon\}$  is a language over  $\Sigma = \{0, 1\}$ . This is an example of a finite language.
- The set of words with odd length over  $\Sigma = \{1\}$ .
- The set of words that contain the same number of 0s and 1s is a language over  $\Sigma = \{0, 1\}$ .

---

<sup>2</sup>Indeed, this is *exactly* how infinite data structures, such as infinite lists, are realised in lazy languages like Haskell.

<sup>3</sup>Given a set  $A$ ,  $\mathcal{P}(A)$  is the *powerset* of  $A$ ; that is, the set of all possible subsets of  $A$ . For example, if  $A = \{a, b\}$ , then  $\mathcal{P}(A) = \{\emptyset, \{a\}, \{b\}, \{a, b\}\}$ . The number of elements  $|A|$  of a set  $A$ , its *cardinality*, and the number of elements in its power set are related by  $|\mathcal{P}(A)| = 2^{|A|}$ . Hence powerset.

- The set of words that contain the same number of 0s and 1s modulo 2 (i.e., both are even or odd) is a language over  $\Sigma = \{0, 1\}$ .
- The set of palindromes using the English alphabet, e.g. words that read the same forwards and backwards like **abba**. This is a language over  $\{a, b, \dots, z\}$ .
- The set of correct Java programs. This is a language over the set of UNICODE “characters” (which correspond to numbers between 0 and  $17 \cdot 2^{16} - 1$ , less some invalid subranges, 1112062 valid encodings in all).
- The set of programs that, if executed on a Windows machine, prints the text “Hello World!” in a window. This is a language over  $\Sigma = \{0, 1\}$ .

Note the distinction between  $\epsilon$ ,  $\emptyset$ , and  $\{\epsilon\}$ !

- $\epsilon$  denotes the empty word, a *sequence* of symbols of length 0.
- $\emptyset$  denotes the empty *set*, a set with no elements.
- $\{\epsilon\}$  is a set with exactly *one* element: the empty word.

In particular, note that  $\epsilon$  is a different type (a sequence) from  $\emptyset$  and  $\{\epsilon\}$  (that are both sets).

An important operation on  $\Sigma^*$  is concatenation. This is denoted by juxtapositioning (or, if you prefer, by an “invisible operator”): given  $u, v \in \Sigma^*$  we can construct a new word  $uv \in \Sigma^*$  simply by concatenating the two words. We can define this operation by primitive recursion:

$$\begin{aligned}\epsilon v &= v \\ (xu)v &= x(uv)\end{aligned}$$

Concatenation is associative and has unit  $\epsilon$ :

$$\begin{aligned}u(vw) &= (uv)w \\ \epsilon u &= u = u\epsilon\end{aligned}$$

where  $u, v, w$  are words. We use exponent notation to denote concatenation of a word with itself. For example,  $u^2 = uu$ ,  $u^3 = uuu$ , and so on. By definition,  $u^1 = u$  and  $u^0 = \epsilon$ , the unit of concatenation. Thus we can simplify repeated concatenation using familiar-looking laws. For example:  $u^1 u^0 u^2 = u^3$ .

Concatenation of words is extended to concatenation of languages by:

$$MN = \{uv \mid u \in M \wedge v \in N\}$$

For example:

$$\begin{aligned}M &= \{\epsilon, a, aa\} \\ N &= \{b, c\} \\ MN &= \{uv \mid u \in \{\epsilon, a, aa\} \wedge v \in \{b, c\}\} \\ &= \{\epsilon b, \epsilon c, ab, ac, aab, aac\} \\ &= \{b, c, ab, ac, aab, aac\}\end{aligned}$$

Some important properties of language concatenation are:

- Concatenation of languages is associative:

$$L(MN) = (LM)N$$

- Concatenation of languages has zero  $\emptyset$ :

$$L\emptyset = \emptyset = \emptyset L$$

- Concatenation of languages has unit  $\{\epsilon\}$ :

$$L\{\epsilon\} = L = \{\epsilon\}L$$

- Concatenation distributes through set union:

$$\begin{aligned} L(M \cup N) &= LM \cup LN \\ (L \cup M)N &= LN \cup MN \end{aligned}$$

Note that concatenation does *not* distribute through intersection! Counterexample. Let  $L = \{\epsilon, a\}$ ,  $M = \{\epsilon\}$ ,  $N = \{a\}$ . Then:

$$\begin{aligned} L(M \cap N) &= L\emptyset = \emptyset \\ LM \cap LN &= \{\epsilon, a\} \cap \{a, aa\} = \{a\} \end{aligned}$$

Exponent notation is used to denote iterated language concatenation:  $L^1 = L$ ,  $L^2 = LL$ ,  $L^3 = LLL$ , and so on. By definition,  $L^0 = \{\epsilon\}$  (for *any* language, including  $\emptyset$ ), which is the unit for language concatenation (just as  $u^0 = \epsilon$  is the unit for concatenation of words).

The Kleene star can also be applied to languages. This intuitively means language concatenation iterated 0 or more times:

$$L^* = \bigcup_{n=0}^{\infty} L^n$$

Note that  $\epsilon \in L^*$  for any language  $L$ , including  $L = \emptyset$ , the empty language. As an example, if  $L = \{a, ab\}$ , then  $L^* = \{\epsilon, a, ab, aab, aba, aaab, aaba, \dots\}$ .

Alternatively (and more abstractly),  $L^*$  can be described as the least language (with respect to  $\subseteq$ ) that contains  $L$  and the empty word,  $\epsilon$ , and is closed under concatenation:

$$u \in L^* \wedge v \in L^* \implies uv \in L^*$$

Note the subtle difference between using the Kleene star on an alphabet  $\Sigma$ , a set of *symbols*, as in  $\Sigma^*$ , and on using the Kleene star on a language  $L \subseteq \Sigma^*$ , a set of *words*. While the result in both cases is a set of words, the types of the arguments to the two variants of the Kleene star operation differ.

## 2.1 Exercises

### Exercise 2.1

Let the alphabet  $\Sigma = \{3, 5, 7, 9\}$ , and let the language  $L = \{w \mid w \in \Sigma^*, 1 \leq |w| \leq 2\}$ . (If  $w$  is a word,  $|w|$  denotes the length of that word. If  $X$  is a finite set, like an alphabet or finite language,  $|X|$  denotes the number of elements in that set, its *cardinality*.) Answer the following questions:

1. Describe  $L$  in plain English.
2. Enumerate all the words in  $L$ .
3. In general, for an arbitrary alphabet  $\Sigma_1$  and  $0 \leq m \leq n$ , how many words are there in the language  $L_1 = \{w \mid w \in \Sigma_1^*, m \leq |w| \leq n\}$ ? That is, write down an expression for  $|L_1|$ .
4. How many words would there be in  $L_1$  if  $\Sigma_1 = \Sigma$ ,  $m = 3$ , and  $n = 7$ ?

### Exercise 2.2

Let the alphabet  $\Sigma = \{a, b, c\}$  and let  $L_1 = \{\epsilon, b, ac\}$  and  $L_2 = \{a, b, ca\}$  be two languages over  $\Sigma$ . Enumerate the words in the following languages, showing your calculations in some detail:

1.  $L_3 = L_1 \cup L_2$
2.  $L_4 = L_1\{\epsilon\}(L_2 \cap L_1)$
3.  $L_5 = L_3\emptyset L_4$

### Exercise 2.3

Let the alphabet  $\Sigma = \{a, b, c\}$  and let  $L_1 = \{\epsilon, b, bb\}$  and  $L_2 = \{a, ab, abc\}$  be two languages over  $\Sigma$ . Enumerate the words in the following languages, showing your calculations in some detail:

1.  $L_3 = L_1 \cap L_2$
2.  $L_4 = (L_2\{\epsilon\}L_1) \cap \Sigma^*$
3.  $L_5 = L_3\emptyset \cap L_4$

### Exercise 2.4

Let the alphabet  $\Sigma = \{a, b, c\}$ . Enumerate the words in

$$L = \{w \mid w \in \{\epsilon, a, b, bc\}^*, |w| \leq 3\}$$

## 3 Finite Automata

Finite automata correspond to a computer with a fixed finite amount of memory<sup>4</sup>. We will introduce *deterministic finite automata* (DFA) first and then move to *nondeterministic finite automata* (NFA). An automaton will accept certain words (sequences of symbols of a given alphabet  $\Sigma$ ) and reject others. The set of accepted words is called the language of the automaton. We will show that the class of languages that are accepted by DFAs and NFAs is the same.

### 3.1 Deterministic finite automata

#### 3.1.1 What is a DFA?

A *deterministic finite automaton* (DFA)  $A = (Q, \Sigma, \delta, q_0, F)$  is given by:

1. A finite set of *states*  $Q$
2. A finite set of input symbols, the *alphabet*,  $\Sigma$
3. A *transition function*  $\delta \in Q \times \Sigma \rightarrow Q$
4. An *initial state*  $q_0 \in Q$
5. A set of *final states*  $F \subseteq Q$

The initial states are sometimes called *start states*, and the final states are sometimes called *accepting states*.

As an example consider the following automaton

$$D = (\{q_0, q_1, q_2\}, \{0, 1\}, \delta_D, q_0, \{q_2\})$$

where

$$\delta_D = \{((q_0, 0), q_1), ((q_0, 1), q_0), ((q_1, 0), q_1), ((q_1, 1), q_2), ((q_2, 0), q_2), ((q_2, 1), q_2)\}$$

if we view a function as a set of argument-result pairs. Alternatively, we can define it case by case:

$$\begin{aligned} \delta_D(q_0, 0) &= q_1 \\ \delta_D(q_0, 1) &= q_0 \\ \delta_D(q_1, 0) &= q_1 \\ \delta_D(q_1, 1) &= q_2 \\ \delta_D(q_2, 0) &= q_2 \\ \delta_D(q_2, 1) &= q_2 \end{aligned}$$

A DFA may be more conveniently represented by a *transition table*. The transition table for the DFA  $D$  is:

$\delta_D$		0	1
$\rightarrow$	$q_0$	$q_1$	$q_0$
	$q_1$	$q_1$	$q_2$
*	$q_2$	$q_2$	$q_2$

---

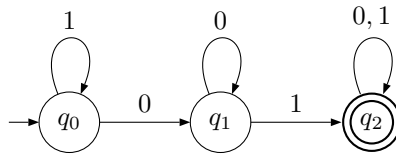
<sup>4</sup>However, that does not mean that finite automata are a good model of general purpose computers. A computer with  $n$  bits of memory has  $2^n$  possible states. That is an absolutely enormous number even for very modest memory sizes, say 1024 bits or more, meaning that describing a computer using finite automata quickly becomes infeasible. We will encounter a better model of computers later, the *Turing Machines*.

A transition table represents the transition function  $\delta$  of a DFA; i.e., the value of  $\delta(q, x)$  is given by the row labelled  $q$  in the column labelled  $x$ . In addition, the initial state is identified by putting an arrow  $\rightarrow$  to the left of it, and all final states are similarly identified by a star  $*$ . The inclusion of this additional information makes a transition table a self-contained representation of a DFA.

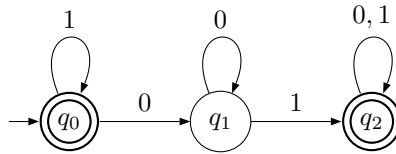
Note that the initial state also can be final (accepting). For example, for a variation  $D'$  of the DFA  $D$  where  $q_0$  also is final:

$\delta_{D'}$		0	1
$\rightarrow *$	$q_0$	$q_1$	$q_0$
	$q_1$	$q_1$	$q_2$
	$*$ $q_2$	$q_2$	$q_2$

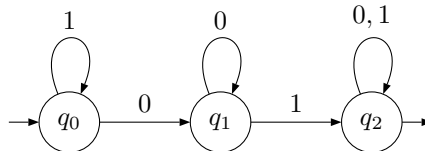
Another way to represent a DFA is through a *transition diagram*. The transition diagram for the DFA  $D$  is:



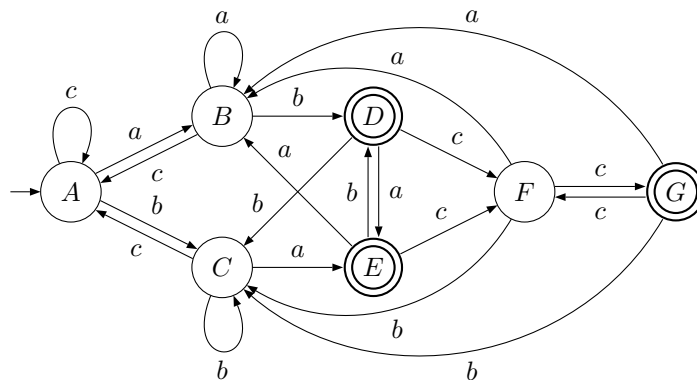
The initial state is identified by an incoming arrow. Final states are drawn with a double outline. If  $\delta(q, x) = q'$  then there is an arrow from state  $q$  to  $q'$  that is labelled  $x$ . For another example, here is the transition diagram for the DFA  $D'$ :



An alternative to the double outline for a final state is to use an outgoing arrow. Using that convention, the transition diagram for the DFA  $D$  is:



Here is an example of a larger DFA over the alphabet  $\Sigma = \{a, b, c\}$  represented by a transition diagram:



The states are named by capital letters this time for a bit of variation:  $Q = \{A, B, C, D, E, F, G\}$ . While it is common to use symbols  $q_i$ ,  $i \in \mathbb{N}$  to name states, we can pick any names we like. Another common choice is to use natural numbers; i.e.,  $Q \subset \mathbb{N} \wedge Q$  is finite.

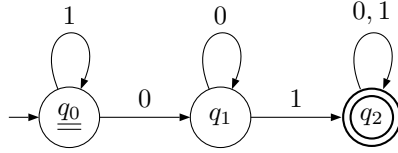
The representation of the above DFA as a transition table is:

$\delta$		$a$	$b$	$c$
$\rightarrow$	$A$	$B$	$C$	$A$
	$B$	$B$	$D$	$A$
	$C$	$E$	$C$	$A$
*	$D$	$E$	$C$	$F$
*	$E$	$B$	$D$	$F$
	$F$	$B$	$C$	$G$
*	$G$	$B$	$C$	$F$

### 3.1.2 The language of a DFA

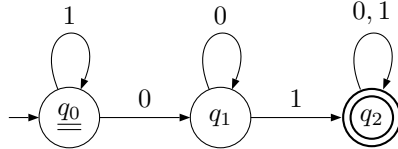
We will now discuss how a DFA accepts or rejects words over its alphabet of input symbols. The set of words accepted by a DFA  $A$  is called the *language*  $L(A)$  of the DFA. Thus, for a DFA  $A$  with alphabet  $\Sigma$ ,  $L(A) \subseteq \Sigma^*$ .

To determine whether a word  $w \in L(A)$ , the machine starts in its initial state. Taking the DFA  $D$  above as an example, it would start in state  $q_0$ . We indicate the state of a DFA by underlining the state name:

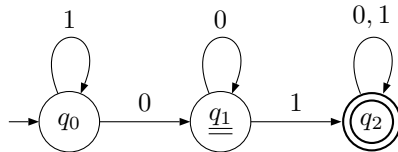


Then, whenever an input symbol is read from  $w$ , the machine transitions to a new state by following the edge labelled with this symbol. Once all symbols in the input word  $w$  have been read, the word is *accepted* if the state is final, meaning  $w \in L(A)$ , otherwise the word is *rejected*, meaning  $w \notin L(A)$ .

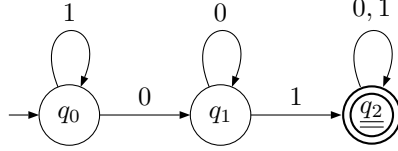
To continue with the example, suppose  $w = 101$ . The machine  $D$  would thus first read 1 and transition to a new state by following the edge labelled 1. As that edge in this case forms a loop back to state  $q_0$ , the machine  $D$  would transition back into state  $q_0$ :



The machine would then read 0 and transition to state  $q_1$  by following the edge labelled 0. We indicate this by moving the mark along that edge to  $q_1$ :



Finally, the machine would read the last 1 in the input word, moving to  $q_2$ :



As  $q_2$  is a final state, the DFA  $D$  accepts the word  $w = 101$ , meaning  $101 \in L(D)$ . In the same way, we can determine that  $0 \notin L(D)$ ,  $110 \notin L(D)$ , but  $011 \in L(D)$ . Verify this. Indeed, a little bit of thought reveals that

$$L(D) = \{w \mid w \text{ contains the substring } 01\}$$

To make the notion of the language of a DFA precise, we now give a formal definition of  $L(A)$ . First we define the *extended transition function*  $\hat{\delta} \in Q \times \Sigma^* \rightarrow Q$ . Intuitively,  $\hat{\delta}(q, w) = q'$  if the machine starting from state  $q$  ends up in state  $q'$  when reading the word  $w$ . Formally,  $\hat{\delta}$  is defined by primitive recursion:

$$\hat{\delta}(q, \epsilon) = q \tag{1}$$

$$\hat{\delta}(q, xw) = \hat{\delta}(\delta(q, x), w) \tag{2}$$

where  $x \in \Sigma$  and  $w \in \Sigma^*$ . Thus,  $xw$  stands for a nonempty word the first symbol of which is  $x$  and the rest of which is  $w$ . For example, if  $xw = 010$ , then  $x = 0$  and  $w = 10$ . Note that  $w$  may be empty; e.g., if  $xw = 0$ , then  $x = 0$  and  $w = \epsilon$ .

As an example, we calculate  $\hat{\delta}_D(q_0, 101) = q_1$ :

$$\begin{aligned}
 \hat{\delta}_D(q_0, 101) &= \hat{\delta}_D(\delta_D(q_0, 1), 01) && \text{by (2)} \\
 &= \hat{\delta}_D(q_0, 01) && \text{because } \delta_D(q_0, 1) = q_0 \\
 &= \hat{\delta}_D(\delta_D(q_0, 0), 1) && \text{by (2)} \\
 &= \hat{\delta}_D(q_1, 1) && \text{because } \delta_D(q_0, 0) = q_1 \\
 &= \hat{\delta}_D(\delta_D(q_1, 1), \epsilon) && \text{by (2)} \\
 &= \hat{\delta}_D(q_2, \epsilon) && \text{because } \delta_D(q_1, 1) = q_2 \\
 &= q_2 && \text{by (1)}
 \end{aligned}$$

Using the extended transition function  $\hat{\delta}$ , we define the language  $L(A)$  of a DFA  $A$  formally:

$$L(A) = \{w \mid \hat{\delta}(q_0, w) \in F\}$$

Returning to our example, we thus have that  $101 \in L(D)$  because  $\hat{\delta}_D(q_0, 101) = q_2$  and  $q_2 \in F_D$ .

## 3.2 Nondeterministic finite automata

### 3.2.1 What is an NFA?

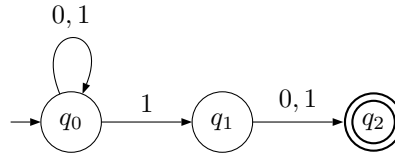
*Nondeterministic finite automata* (NFA) have transition functions that map a given state and an input symbol to zero or more successor states. We can think of this as the machine having a “choice” whenever there are two or more possible transitions from a state on an input symbol. In this presentation, we will further



allow an NFA to have more than one initial state<sup>5</sup>. Again, we can think of this as the machine having a “choice” of where to start. An NFA accepts a word  $w$  if there is at least one *possible* way to get from one of the initial states to one of the final states along edges labelled with the symbols of  $w$  in order.

It is important to note that although an NFA has a nondeterministic transition function, it can always be determined whether or not a word belongs to its language. Indeed, we shall see that every NFA can be translated into an DFA that accepts the same language.

Here is an example of an NFA  $C$  that accepts all words over  $\Sigma = \{0, 1\}$  such that the symbol before the last is 1:



A *nondeterministic finite automaton* (NFA)  $A = (Q, \Sigma, \delta, S, F)$  is given by:

1. A finite set of *states*  $Q$ ,
2. A finite set of input symbols, the *alphabet*,  $\Sigma$ ,
3. A *transition function*  $\delta \in Q \times \Sigma \rightarrow \mathcal{P}(Q)$ ,
4. A set of *initial states*  $S \subseteq Q$ ,
5. A set of *final* (or *accepting*) states  $F \subseteq Q$ .

Thus, in contrast to a DFA, an NFA may have many initial states, not just one, and its transition function maps a state and an input symbol to a set of possible successor states, not just a single state. As an example we have that

$$C = (\{q_0, q_1, q_2\}, \{0, 1\}, \delta_C, \{q_0\}, \{q_2\})$$

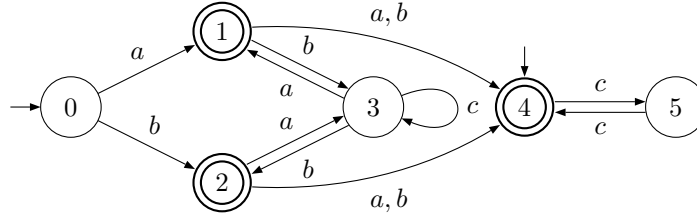
where  $\delta_C$  is given by

$\delta_C$		0	1
$\rightarrow$	$q_0$	$\{q_0\}$	$\{q_0, q_1\}$
	$q_1$	$\{q_2\}$	$\{q_2\}$
*	$q_2$	$\emptyset$	$\emptyset$

Note that the entries in the table are *sets* of states, and that these sets may be empty ( $\emptyset$ ), here exemplified by the entries for state  $q_2$ . Again, the (in this case only) initial state has been marked with  $\rightarrow$  and the (in this case only) final state marked with  $*$  to make this a self-contained representation of the NFA.

Here is another example of an NFA, this time over the alphabet  $\Sigma = \{a, b, c\}$  and with states  $Q = \{0, 1, 2, 3, 4, 5\} \subset \mathbb{N}$ :

<sup>5</sup>Note that we diverge slightly from the definition in the book [HMu01], which uses a single initial state instead of a set of initial states. Permitting more than one initial state allows us to avoid introducing  $\epsilon$ -NFAs (see [HMu01], section 2.5).



The transition table for this NFA is:

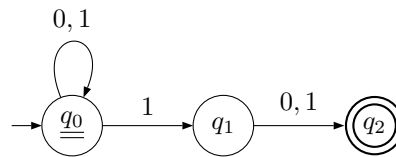
$\delta$		$a$	$b$	$c$
$\rightarrow$	0	{1}	{2}	$\emptyset$
*	1	{4}	{3, 4}	$\emptyset$
*	2	{3, 4}	{4}	$\emptyset$
	3	{1}	{2}	{3}
$\rightarrow$	4	$\emptyset$	$\emptyset$	{5}
	5	$\emptyset$	$\emptyset$	{4}

Note that this NFA has multiple initial states, multiple final states, one initial state that also is final, and that there in some cases are no possible successor states and in other cases more than one.

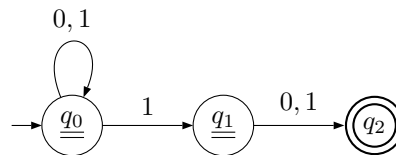
### 3.2.2 The language accepted by an NFA

To see whether a word  $w \in \Sigma^*$  is accepted by an NFA  $A$ , we have to consider *all* possible states the machine could be in after having read a sequence of input symbols. Initially, an NFA can be in any of its initial states. Each time an input symbol is read, all successor states on the read symbol for each current possible state become the new possible states. After having read a complete word  $w$ , if at least one of the possible states is final (accepting), then that word is accepted, meaning  $w \in L(A)$ , otherwise it is rejected, meaning  $w \notin L(A)$ .

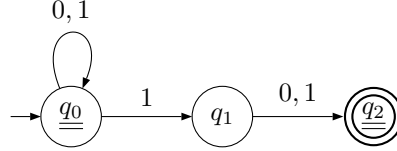
We will illustrate by showing how the NFA  $C$  rejects the word 100. We will again mark the current states of the NFA by underlining the state names, but this time there may be more than one marked state at once. Initially, as  $q_0$  is the only initial state, we would have:



Each time when we read a symbol we look at all the marked states. We remove the old markers and put markers at all the states that are *reachable* via an arrow marked with the current input symbol. This may include one or more states that were marked previously. It may also be the case that no states are reachable, in which case all marks are removed and the word rejected (as it no longer is possible to reach any final states). In our example, after reading 1, there would be two marked states as there are two arrows from  $q_0$  labelled 1:

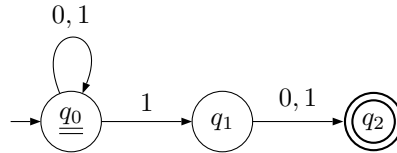


After reading 0, the next symbol in the word 100, there would still be two marked states as the machine on input 0 can reach  $q_0$  from  $q_0$  and  $q_2$  from  $q_1$ :



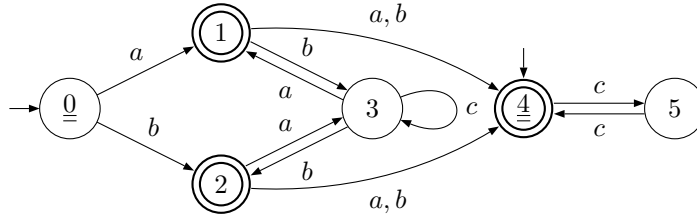
Note that one of the marked states is a final (accepting) state, meaning the word read so far (10) is accepted by the NFA.

However, there is one symbol left in our example word 100, and after having read the final 0, the final state would no longer be marked because it cannot be reached from any of the marked states:

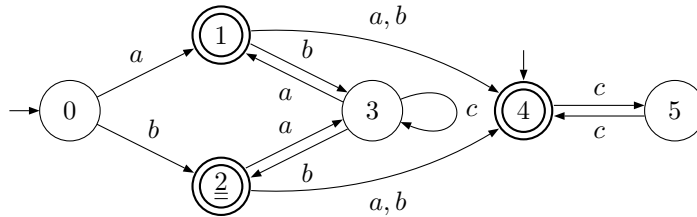


The NFA  $C$  thus rejects the word 100.

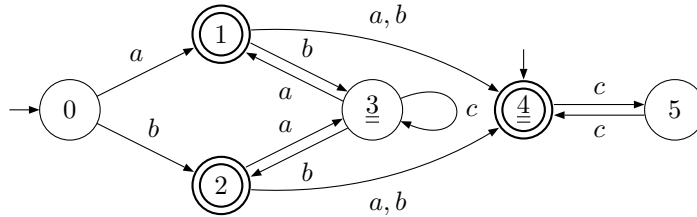
For another example, consider the NFA at the end of section 3.2.1. Convince yourself that you understand how this NFA accepts the words  $\epsilon$ ,  $abcc$ ,  $abcca$ , and rejects  $abccaac$ . We illustrate by tracing its operation on the word  $bacac$ . We start by marking all initial states. Then it is just a matter of systematically exploring all possibilities:



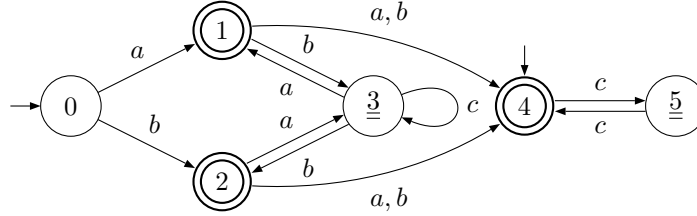
After reading  $b$ :



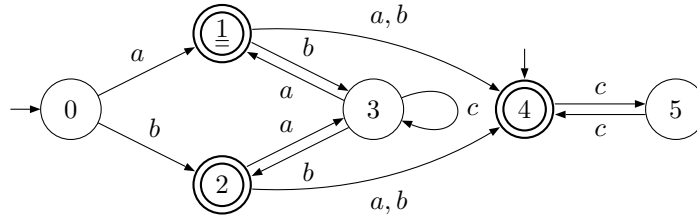
After reading  $a$ :



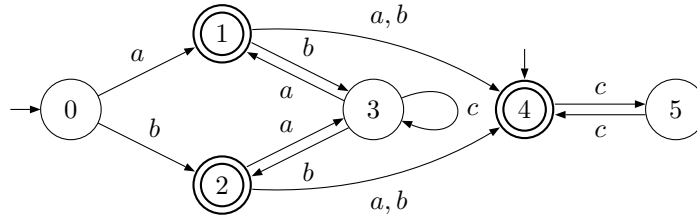
After reading  $c$ :



After reading  $a$ :



After reading  $c$ :



The machine thus rejects  $bacac$  as no final state is marked. In fact, as there are *no* marked states left at all, this shows that this NFA will reject *all* words that start  $bacac$ . Can you find other such prefixes?

To define the *extended transition function*  $\hat{\delta}$  for NFAs we use a generalisation of the union operation  $\cup$  on sets over a (finite) set of sets:

$$\bigcup \{A_1, A_2, \dots, A_n\} = A_1 \cup A_2 \cup \dots \cup A_n$$

In the special cases of the empty set of sets and a one element set of sets:

$$\bigcup \emptyset = \emptyset \quad \bigcup \{A\} = A$$

As an example

$$\bigcup \{\{1\}, \{2, 3\}, \{1, 3\}\} = \{1\} \cup \{2, 3\} \cup \{1, 3\} = \{1, 2, 3\}$$

Alternatively, we can define  $\bigcup$  by comprehension, which also extends the operation to infinite sets of sets (although we don't need this here):

$$\bigcup B = \{x \mid \exists A \in B. x \in A\}$$

We define  $\hat{\delta} \in \mathcal{P}(Q) \times \Sigma^* \rightarrow \mathcal{P}(Q)$  such that  $\hat{\delta}(P, w)$  is the set of states that are reachable from one of the states in  $P$  on the word  $w$ :

$$\hat{\delta}(P, \epsilon) = P \quad (3)$$

$$\hat{\delta}(P, xw) = \hat{\delta}(\bigcup \{\delta(q, x) \mid q \in P\}, w) \quad (4)$$

where  $x \in \Sigma$  and  $w \in \Sigma^*$ . Intuitively, if  $P$  are the possible states, then  $\hat{\delta}(P, w)$  are the possible states after having read a word  $w$ .

To illustrate, we calculate  $\hat{\delta}_C(q_0, 100)$ :

$$\begin{aligned} \hat{\delta}_C(\{q_0\}, 100) &= \hat{\delta}_C(\bigcup \{\delta_C(q, 1) \mid q \in \{q_0\}\}, 00) && \text{by (4)} \\ &= \hat{\delta}_C(\delta_C(q_0, 1), 00) \\ &= \hat{\delta}_C(\{q_0, q_1\}, 00) \\ &= \hat{\delta}_C(\bigcup \{\delta_C(q, 0) \mid q \in \{q_0, q_1\}\}, 0) && \text{by (4)} \\ &= \hat{\delta}_C(\delta_C(q_0, 0) \cup \delta_C(q_1, 0), 0) \\ &= \hat{\delta}_C(\{q_0\} \cup \{q_2\}, 0) \\ &= \hat{\delta}_C(\{q_0, q_2\}, 0) \\ &= \hat{\delta}_C(\bigcup \{\delta_C(q, 0) \mid q \in \{q_0, q_2\}\}, \epsilon) && \text{by (4)} \\ &= \hat{\delta}_C(\delta_C(q_0, 0) \cup \delta_C(q_2, 0), \epsilon) \\ &= \hat{\delta}_C(\{q_0\} \cup \emptyset, \epsilon) \\ &= \{q_0\} && \text{by (3)} \end{aligned}$$

Of course, we already knew this from the worked example above illustrating how the NFA  $C$  rejects 100. Make sure you see how the marked states after each step coincides with the set of possible states in the calculation.

The language of an NFA can now be defined using  $\hat{\delta}$ :

$$L(A) = \{w \mid \hat{\delta}(S, w) \cap F \neq \emptyset\}$$

Thus,  $100 \notin L(C)$  because

$$\hat{\delta}_C(S_C, 100) \cap F_C = \hat{\delta}_C(\{q_0\}, 100) \cap \{q_2\} = \{q_0\} \cap \{q_2\} = \emptyset$$

### 3.2.3 The subset construction

DFAs can be viewed as a special case of NFAs; i.e., those for which there is precisely one start state ( $S = \{q_0\}$ ) and for which the transition function always returns singleton (one-element) sets ( $\delta(q, x) = \{q'\}$  for all  $q \in Q$  and  $x \in \Sigma$ ).

The opposite is also true, however: NFAs are really just DFAs “in disguise”. We show this by for a given NFA systematically constructing an *equivalent* DFA; i.e., a DFA that accepts the same language as the given NFA. NFAs are thus no more powerful than DFAs; i.e., NFAs cannot describe more languages than DFAs. However, in some cases, NFAs need a lot fewer states than the corresponding DFA, and they may be easier to construct in the first place.

*The subset construction:* Given an NFA  $A = (Q, \Sigma, \delta, S, F)$  we construct the equivalent DFA:

$$D(A) = (\mathcal{P}(Q), \Sigma, \delta_{D(A)}, S, F_{D(A)})$$

where

$$\delta_{D(A)}(P, x) = \bigcup \{\delta(q, x) \mid q \in P\} \quad (5)$$

$$F_{D(A)} = \{P \in \mathcal{P}(Q) \mid P \cap F \neq \emptyset\} \quad (6)$$

The basic idea of this construction is to define a DFA whose states are *sets of NFA states*. A set of possible NFA states thus becomes a single DFA state. The DFA transition function is given by considering all reachable NFA states for each of the current possible NFA states for each input symbol. The resulting set of possible NFA states is again just a single DFA state. A DFA state is final if that set that contains at least one final NFA state.

As an example, let us construct a DFA  $D(C)$  equivalent to  $C$  above:

$$D(C) = (\mathcal{P}(\{q_0, q_1, q_2\}), \{0, 1\}, \delta_{D(C)}, \{q_0\}, F_{D(C)})$$

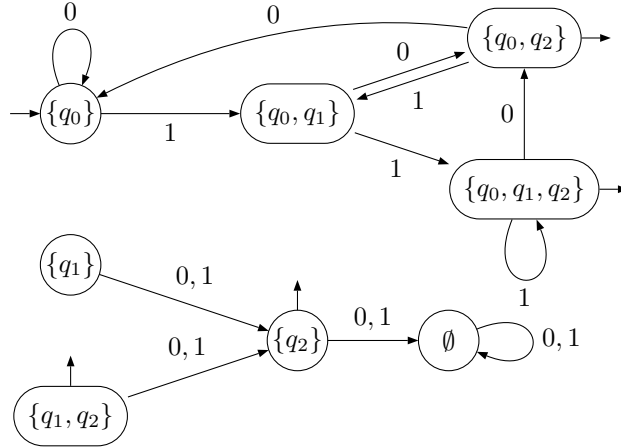
where  $\delta_{D(C)}$  is given by:

$\delta_{D(C)}$		0	1
	$\emptyset$	$\emptyset$	$\emptyset$
$\rightarrow$	$\{q_0\}$	$\{q_0\}$	$\{q_0, q_1\}$
	$\{q_1\}$	$\{q_2\}$	$\{q_2\}$
*	$\{q_2\}$	$\emptyset$	$\emptyset$
	$\{q_0, q_1\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$
*	$\{q_0, q_2\}$	$\{q_0\}$	$\{q_0, q_1\}$
*	$\{q_1, q_2\}$	$\{q_2\}$	$\{q_2\}$
*	$\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$

and  $F_{D(C)}$  (all the states marked with \* above) by:

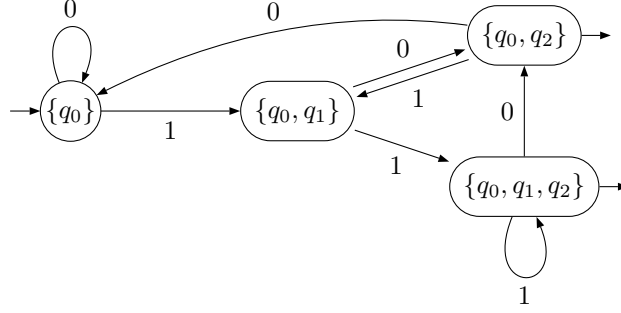
$$F_{D(C)} = \{\{q_2\}, \{q_0, q_2\}, \{q_1, q_2\}, \{q_0, q_1, q_2\}\}$$

The transition diagram is:



Accepting states have been marked by outgoing arrows.

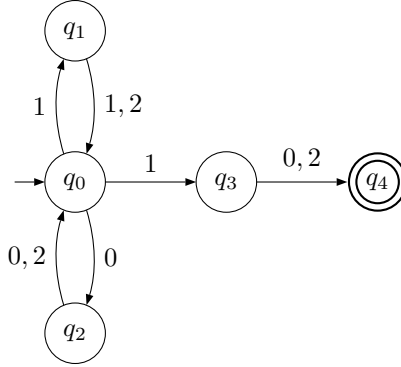
Note that some of the states ( $\emptyset, \{q_1\}, \{q_2\}, \{q_1, q_2\}$ ) cannot be reached from the initial state. This means that they can be omitted without changing the language. We thus obtain the following automaton:



It is possible to avoid having to perform calculations for states that cannot be reached by carrying out the subset construction in a “demand driven” way. The idea is to start from the initial DFA state, which is just the set of initial NFA states  $S$ , and then only consider the DFA states (subsets of NFA states) that appear during the course of the calculations. We illustrate this approach by an example. Consider the following NFA  $N$ :

$$N = (Q_N = \{q_0, q_1, q_2, q_3, q_4\}, \Sigma_N = \{0, 1, 2\}, \delta_N, S_N = \{q_0\}, F_N = \{q_4\})$$

where  $\delta_N$  is given by the transition diagram:



Note that  $N$  has 5 states which means that the DFA  $D(N)$  has  $|\mathcal{P}(Q_N)| = 2^5 = 32$  states. However, as we will see, only a handful of those 32 states can actually be reached from the initial state  $S_N$  of  $D(N)$ . We would thus waste quite a bit of effort if we were to tabulate all of them.

We start from  $S_N = \{q_0\}$ , the set of start states of  $N$ , and we compute  $\bigcup\{\delta(q, x) \mid q \in S_N\}$  for each  $x \in \Sigma_N$  (equation (5)). In this case we get:

$\delta_{D(N)}$	0	1	2
$\rightarrow \{q_0\}$	$\{q_2\}$	$\{q_1, q_3\}$	$\emptyset$

Whenever we encounter a state  $P \subseteq Q$  of  $D(N)$  that has not been considered before, we add  $P$  to the table, marking any final states as such. In this case, three new DFA states emerge ( $\{q_2\}$ ,  $\{q_1, q_3\}$ , and  $\emptyset$ ), none of which is final:

$\delta_{D(N)}$	0	1	2
$\rightarrow \{q_0\}$	$\{q_2\}$	$\{q_1, q_3\}$	$\emptyset$
$\{q_2\}$			
$\{q_1, q_3\}$			
$\emptyset$			

We then proceed to tabulate  $\delta_{D(N)}$  for each of the new states for each  $x \in \Sigma$ , adding any further new states to the table:

$\delta_{D(N)}$		0	1	2
$\rightarrow \{q_0\}$		$\{q_2\}$	$\{q_1, q_3\}$	$\emptyset$
$\{q_2\}$		$\{q_0\}$	$\emptyset$	$\{q_0\}$
$\{q_1, q_3\}$	$\emptyset \cup \{q_4\} = \{q_4\}$	$\{q_0\} \cup \emptyset = \{q_0\}$	$\{q_0\} \cup \{q_4\} = \{q_0, q_4\}$	
$\emptyset$		$\emptyset$	$\emptyset$	$\emptyset$

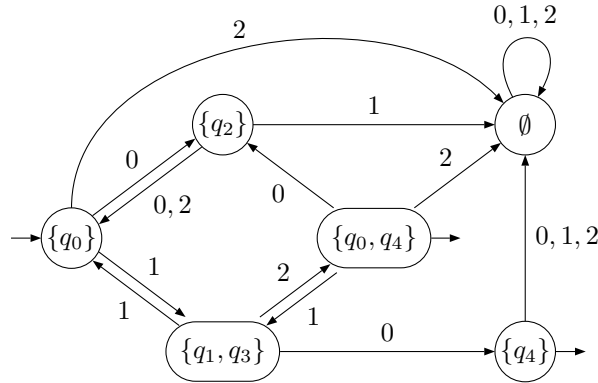
Here, two new states emerge ( $\{q_4\}$  and  $\{q_0, q_4\}$ ), both final (because  $\{q_4\} \cap F_N \neq \emptyset$  and  $\{q_0, q_4\} \cap F_N \neq \emptyset$ ):

$\delta_{D(N)}$		0	1	2
$\rightarrow \{q_0\}$		$\{q_2\}$	$\{q_1, q_3\}$	$\emptyset$
$\{q_2\}$		$\{q_0\}$	$\emptyset$	$\{q_0\}$
$\{q_1, q_3\}$	$\emptyset \cup \{q_4\} = \{q_4\}$	$\{q_0\} \cup \emptyset = \{q_0\}$	$\{q_0\} \cup \{q_4\} = \{q_0, q_4\}$	
$\emptyset$		$\emptyset$	$\emptyset$	$\emptyset$
* $\{q_4\}$				
* $\{q_0, q_4\}$				

This process is repeated until no new states emerges. Tabulating for the last two new states reveals that no further states emerge in this case and we are thus done, having only had to tabulate for 6 reachable out of the 32 DFA states:

$\delta_{D(N)}$		0	1	2
$\rightarrow \{q_0\}$		$\{q_2\}$	$\{q_1, q_3\}$	$\emptyset$
$\{q_2\}$		$\{q_0\}$	$\emptyset$	$\{q_0\}$
$\{q_1, q_3\}$	$\emptyset \cup \{q_4\} = \{q_4\}$	$\{q_0\} \cup \emptyset = \{q_0\}$	$\{q_0\} \cup \{q_4\} = \{q_0, q_4\}$	
$\emptyset$		$\emptyset$	$\emptyset$	$\emptyset$
* $\{q_4\}$		$\emptyset$	$\emptyset$	$\emptyset$
* $\{q_0, q_4\}$	$\{q_2\} \cup \emptyset = \{q_2\}$	$\{q_1, q_3\} \cup \emptyset = \{q_1, q_3\}$	$\emptyset \cup \emptyset = \emptyset$	

After double checking that we have not forgotten to mark any final states, we can draw the transition diagram for  $D(N)$ :



Accepting states have been marked by outgoing arrows.



### 3.2.4 Correctness of the subset construction

We still have to convince ourselves that the subset construction actually works; i.e., that for a given NFA  $A$  it really is the case that  $L(A) = L(D(A))$ . We start by proving the following lemma, which says that the extended transition functions coincide:

**Lemma 3.1**

$$\hat{\delta}_{D(A)}(P, w) = \hat{\delta}_A(P, w)$$

The result of both functions is a set of states of the NFA  $A$ : for the left-hand side because the extended transition function on NFAs returns a set of states, and for the right-hand side because the states of  $D(A)$  are sets of states of  $A$ .

**Proof.** We show this by *induction* over the length of the word  $w$ ,  $|w|$ .

$|w| = 0$  Then  $w = \epsilon$  and we have

$$\begin{aligned} \hat{\delta}_{D(A)}(P, \epsilon) &= P && \text{by (1)} \\ &= \hat{\delta}_A(P, \epsilon) && \text{by (3)} \end{aligned}$$

$|w| = n + 1$  Then  $w = xv$  with  $|v| = n$ .

$$\begin{aligned} \hat{\delta}_{D(A)}(P, xv) &= \hat{\delta}_{D(A)}(\delta_{D(A)}(P, x), v) && \text{by (2)} \\ &= \hat{\delta}_A(\delta_{D(A)}(P, x), v) && \text{ind.hyp.} \\ &= \hat{\delta}_A(\bigcup\{\delta_A(q, x) \mid q \in P\}, v) && \text{by (5)} \\ &= \hat{\delta}_A(P, xv) && \text{by (4)} \end{aligned}$$

□

We can now use the lemma to show

**Theorem 3.2**

$$L(A) = L(D(A))$$

**Proof.**

$$\begin{aligned} &w \in L(A) \\ \iff &\text{Definition of } L(A) \text{ for NFAs} \\ &\hat{\delta}_A(S, w) \cap F \neq \emptyset \\ \iff &\text{Lemma 3.1} \\ &\hat{\delta}_{D(A)}(S, w) \cap F \neq \emptyset \\ \iff &\text{Definition of } F_{D(A)} \\ &\hat{\delta}_{D(A)}(S, w) \in F_{D(A)} \\ \iff &\text{Definition of } L(A) \text{ for DFAs} \\ &w \in L_{D(A)} \end{aligned}$$

□

**Corollary 3.3** *NFAs and DFAs recognise the same class of languages.*

**Proof.** We have noticed that DFAs are just a special case of NFAs. On the other hand the subset construction introduced above shows that for every NFA we can find a DFA that recognises the same language. □

### 3.3 Exercises

#### Exercise 3.1

Let the alphabet  $\Sigma_A = \{a, b\}$  and consider the following DFA  $A$ :

$$\begin{aligned} A &= (Q_A = \{0, 1, 2, 3\}, \Sigma_A, \delta_A, q_0 = 0, F_A = \{1, 2\}) \\ \delta_A &= \{((0, a), 1), ((0, b), 2), ((1, a), 0), ((1, b), 3), ((2, a), 3), ((2, b), 0), \\ &\quad ((3, a), 2), ((3, b), 1)\} \end{aligned}$$

(Here tuple notation is used to define the mapping of the transition function  $\delta_A$ ; thus  $\delta_A(0, a) = 1$ ,  $\delta_A(0, b) = 2$ , etc.) For the DFA  $A$ :

1. Draw its transition diagram.
2. Determine which of the following words belong to  $L(A)$ :
  1.  $\epsilon$
  2.  $b$
  3.  $abaab$
  4.  $bababbba$
3. Explicitly calculate  $\hat{\delta}_A(0, abba)$ .
4. Describe the language that the automaton recognises in English.

#### Exercise 3.2

Construct a DFA  $B$  over  $\Sigma_B = \{a, b, c, d\}$  accepting all words where the number of  $a$ 's is a multiple of 3. E.g.  $abdaca \in L(B)$  (3  $a$ 's), but  $ddaabaa \notin L(B)$  (4  $a$ 's, 4 is not a multiple of 3). Explain your construction. In particular, explain why you chose to have the number of states you did, and explain the purpose (or "meaning") of each state.

#### Exercise 3.3

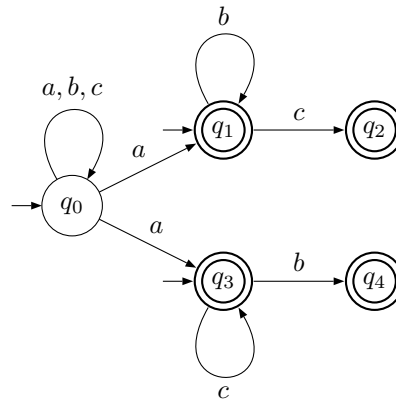
For the alphabet  $\Sigma_C = \{a, b, c\}$ , construct a DFA  $C$  that recognises all words where the number of  $a$ 's is odd and the number of  $b$ 's is divisible by 3. (There may thus be any number of  $c$ 's.) For example,  $a \in L(C)$  (odd number of  $a$ 's, the number of  $b$ 's is 0 which is divisible by 3),  $cbcbabc \in L(C)$  (odd number of  $a$ 's, the number of  $b$ 's is 3 which is divisible by 3), but  $\epsilon \notin L(C)$  (even number of  $a$ 's). Give a brief explanation of your construction, that clearly conveys the key ideas, and give the transition diagram for your DFA as the final answer.

#### Exercise 3.4

For the alphabet  $\Sigma_D = \{0, 1, 2, 3\}$ , construct a DFA  $D$  that precisely recognizes the words for which the arithmetic sum of the constituent symbols is divisible by 5. For example,  $\epsilon \in L(D)$  (there are no symbols in the empty string, the sum is thus 0 which is divisible by 5),  $0 \in L(D)$  (the sum is again 0), and  $23131 \in L(D)$  ( $2 + 3 + 1 + 3 + 1 = 10$  which is divisible by 5), but  $133 \notin L(D)$  ( $1 + 3 + 3 = 7$  which is not divisible by 5). Explain your construction.

### Exercise 3.5

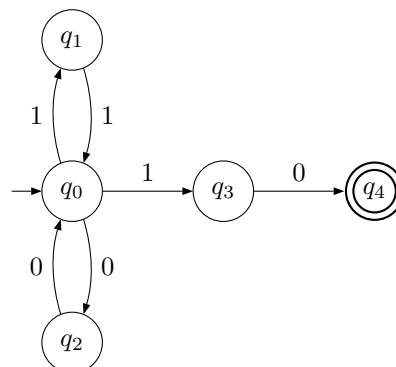
Consider the following NFA  $A$  over  $\Sigma_A = \{a, b, c\}$ :



- Which of the following words are accepted by  $A$  and which are not?
  - $\epsilon$
  - $aaa$
  - $bbc$
  - $cbc$
  - $abcacb$
- Construct a DFA  $D(A)$  equivalent to  $A$  using the “subset construction”. Clearly show each step of your calculations in a transition table.  
*Hint:* Some of the 32 states (i.e., the  $2^{|Q_A|} = 2^5 = 32$  possible subsets of  $Q_A$ ) that would arise by applying the subset construction blindly to  $A$  may be unreachable. You can therefore adopt a strategy where you only consider states reachable from the initial state,  $S_A$ .
- Draw the transition diagram for  $D(A)$ , ignoring unreachable states.

### Exercise 3.6

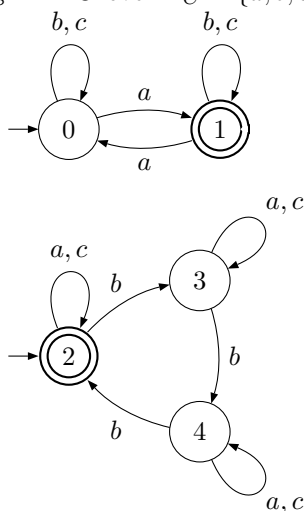
Consider the following NFA  $B$  over  $\Sigma_B = \{0, 1\}$ :



1. Construct a DFA  $D(B)$  equivalent to  $B$  using the “subset construction” and draw the transition diagram for  $D(B)$ , ignoring unreachable states. Clearly show each step of your calculations, e.g. in a transition table.
2. Carry out a *sanity check* on your resulting DFA  $D(B)$  as follows.
  - (a) Give two words over  $\Sigma_B$  that are accepted by the NFA  $B$  and two that are not. At least two of those should be four symbols long or longer.
  - (b) Check that the DFA  $D(B)$  accepts the first two words and rejects the other two, exactly like the NFA  $B$ . Justify your answer by listing the sequence of states the DFA  $D(B)$  goes through for each word, and stating whether or not the last state of that sequence is accepting.

### Exercise 3.7

Consider the following NFA  $C$  over  $\Sigma_C = \{a, b, c\}$ :



1. Which of the following words are accepted by  $C$  and which are not?
  - (a)  $\epsilon$
  - (b)  $aa$
  - (c)  $bb$
  - (d)  $abcabc$
  - (e)  $abcbca$
2. Describe the language accepted by  $C$  in English.
3. Construct a DFA  $D(C)$  equivalent to  $C$  using the “subset construction”. Clearly show each step of your calculations in a transition table. Indicate which DFA state that is initial and which DFA states that are accepting. Only consider states reachable from the initial state of the resulting DFA.
4. Draw the transition diagram for  $D(C)$ .

## 4 Regular Expressions

To recapitulate, given an alphabet  $\Sigma$ , a language is a set of words  $L \subseteq \Sigma^*$ . So far, we have described languages either using set theory (explicit enumeration or set comprehensions) or through finite automata. The key benefits of using automata is that they can describe infinite languages (unlike enumeration) and that they directly give a mechanical way to determine language membership (unlike comprehensions). However, from an automaton, it is not usually immediately obvious what the language of that automaton is, and conversely, given a high-level description of a language, it is often not obvious if it is possible to describe the language using a finite automaton.

This section introduces *regular expressions*: a concise and much more direct way to describe languages. Moreover, a regular expression can mechanically be translated into a finite automaton that accept precisely the language described. This opens up for many practical applications as languages can both be described and recognised with ease. In fact, the opposite is also true: given a finite automaton, it is possible to translate that into an equivalent regular expression. Finite automata and regular expressions are thus *interconvertible*, meaning that they describe the exact *same* class of languages: the *regular languages* or, according to the Chomsky hierarchy, type 3 languages (section 1.1).

One application of regular expressions is to define patterns in programs such as **grep**. Given a regular expression and a sequence of text lines as input, **grep** returns those lines that match the regular expression, where matching means that the line contains a substring that is in the language denoted by the regular expression. The syntax used by **grep** for regular expressions is slightly different from the one used here, and **grep** further supports some convenient abbreviations. However, the underlying theory is exactly the same.

Other applications for regular expressions include defining the *lexical* syntax of programming languages; i.e., what basic symbols, or *tokens*, such as identifiers, keywords, numeric literals look like, as well other lexical aspects such as white space and comments. The context-free syntax (see section 7) of a programming language is then defined in terms of the tokens; i.e., the tokens effectively constitute the alphabet of the language.

In fact, regular expression matching has so many applications that many programming languages provide support for this capability, either built-in or via libraries. Examples include Perl, PHP, Python, and Java. In the past, some of those implementations were a bit naive as the regular expressions were not compiled into finite automata. As a result, matching could be *very* slow, as explained in the paper *Regular Expression Matching Can Be Simple And Fast (but is slow in Java, Perl, PHP, Python, Ruby, ...)* [Cox07]. This paper is a very good read, and once you have read these lecture notes up to and including the present section, you will be able to appreciate it fully.

### 4.1 What are regular expressions?

Given an alphabet  $\Sigma$  (e.g.,  $\Sigma = \{a, b, c, \dots, z\}$ ), the *syntax* (i.e., *form*) of regular expressions over  $\Sigma$  is defined inductively as follows:

1.  $\emptyset$  is a regular expression.
2.  $\epsilon$  is a regular expression.

3. For each  $x \in \Sigma$ ,  $\mathbf{x}$  is a regular expression<sup>6</sup>.
4. If  $E$  and  $F$  are regular expressions then  $E + F$  is a regular expression.
5. If  $E$  and  $F$  are regular expressions then  $EF$  (juxtapositioning; just one after the other) is a regular expression.
6. If  $E$  is a regular expression then  $E^*$  is a regular expression.
7. If  $E$  is a regular expression then  $(E)$  is a regular expression<sup>7</sup>.

These are all regular expressions.

To illustrate, here are some examples of regular expressions:

- $\epsilon$
- **hallo**
- **hallo + hello**
- **h(a + e)llo**
- **a\*b\***
- **( $\epsilon + \mathbf{b}$ )(ab)\*( $\epsilon + \mathbf{a}$ )**

As in arithmetic, there are conventions for reading regular expressions:

- $*$  binds stronger than juxtapositioning and  $+$ . For example, **ab\*** is read as **a(b\*)**. Parentheses must be used to enforce the reading **(ab)\***.
- Juxtapositioning binds stronger than  $+$ . For example, **ab + cd** is read as **(ab) + (cd)**. Parentheses must be used to enforce the reading **a(b + c)d**.

## 4.2 The meaning of regular expressions

In the previous section, we defined the *syntax* of regular expressions, their *form*. We now proceed to define the *semantics* of regular expressions; i.e., what they *mean*, what language a regular expression denotes.

To answer this question, first recall the definition of concatenation of concatenation of languages from section 2:

$$L_1 L_2 = \{uv \mid u \in L_1 \wedge v \in L_2\}$$

We further recall the *Kleene star* operation from the same section (2):

$$L^* = \bigcup_{n=0}^{\infty} L^n$$

To each regular expression  $E$  over  $\Sigma$  we assign a language  $L(E) \subseteq \Sigma^*$  as its meaning or semantics. We do this by *induction over the definition of the syntax*:

---

<sup>6</sup>Note that the regular expression here is typeset in boldface, like **a**, to distinguish it from the corresponding *symbol*, like *a*, typeset in a type-writer font in this and the next section (and on occasion later on as well). Underlining is sometimes used as an alternative to boldface.

<sup>7</sup>The parentheses have been typeset in boldface to emphasise that they are part of the syntax of the regular expression.

1.  $L(\emptyset) = \emptyset$
2.  $L(\epsilon) = \{\epsilon\}$
3.  $L(\mathbf{x}) = \{x\}$  where  $x \in \Sigma$ .
4.  $L(E + F) = L(E) \cup L(F)$
5.  $L(EF) = L(E)L(F)$
6.  $L(E^*) = L(E)^*$
7.  $L((E)) = L(E)$

Subtle points: In (1), the symbol  $\emptyset$  is used both as a regular expression and as the empty set (empty language). Similarly,  $\epsilon$  in (2) is used in two ways: as a regular expression and as the empty word. In (3), the regular expression is typeset in boldface to distinguish it from the corresponding symbol. In (6), the  $*$ -operator is used both to construct a regular expression (part of the syntax) and as an operation on languages. In (7), the inner parentheses on the left-hand side, typeset in boldface, are part of the syntax of regular expressions.

Let us now calculate the meaning of each of the regular expression examples from the previous section; i.e., the language denoted in each case:

- $\epsilon$ :

By (2):

$$L(\epsilon) = \{\epsilon\}$$

- **hallo**:

Consider  $L(\mathbf{ha})$ . By (3):

$$\begin{aligned} L(\mathbf{h}) &= \{\mathbf{h}\} \\ L(\mathbf{a}) &= \{\mathbf{a}\} \end{aligned}$$

Hence, by (5) and language concatenation (section 2):

$$\begin{aligned} L(\mathbf{ha}) &= L(\mathbf{h})L(\mathbf{a}) \\ &= \{uv \mid u \in L(\mathbf{h}) \wedge v \in L(\mathbf{a})\} \\ &= \{uv \mid u \in \{\mathbf{h}\} \wedge v \in \{\mathbf{a}\}\} \\ &= \{\mathbf{ha}\} \end{aligned}$$

Continuing the same reasoning we obtain:

$$L(\mathbf{hallo}) = \{\mathbf{hallo}\}$$

- **hallo + hello**:

From above we know  $L(\mathbf{hallo}) = \{\mathbf{hallo}\}$  and  $L(\mathbf{hello}) = \{\mathbf{hello}\}$ . By (4) we then get:

$$\begin{aligned} L(\mathbf{hallo} + \mathbf{hello}) &= \{\mathbf{hallo}\} \cup \{\mathbf{hello}\} \\ &= \{\mathbf{hallo}, \mathbf{hello}\} \end{aligned}$$

- $\mathbf{h(a + e)llo}$ :

By (3) and (4) we know  $L(\mathbf{a + e}) = \{\mathbf{a, e}\}$ . Thus, using (5) and language concatenation, we obtain:

$$\begin{aligned}
 L(\mathbf{h(a + e)llo}) &= L(\mathbf{h})L(\mathbf{a + e})L(\mathbf{llo}) \\
 &= \{uvw \mid u \in L(\mathbf{h}) \wedge v \in L(\mathbf{a + e}) \wedge w \in L(\mathbf{llo})\} \\
 &= \{uvw \mid u \in \{\mathbf{h}\} \wedge v \in \{\mathbf{a, e}\} \wedge w \in \{\mathbf{llo}\}\} \\
 &= \{\mathbf{hallo, hello}\}
 \end{aligned}$$

- $\mathbf{a^*b^*}$ :

By (6):

$$\begin{aligned}
 L(\mathbf{a^*}) &= L(\mathbf{a})^* \\
 &= \{\mathbf{a}\}^* \\
 &= \bigcup_{n=0}^{\infty} \{\mathbf{a}\}^n \\
 &= \bigcup_{n=0}^{\infty} \{w_1 w_2 \dots w_n \mid 1 \leq i \leq n, w_i \in \{\mathbf{a}\}\} \\
 &= \bigcup_{n=0}^{\infty} \{\mathbf{a}^n\} \\
 &= \{\mathbf{a}^n \mid n \in \mathbb{N}\}
 \end{aligned}$$

Using (5) and language concatenation, this allows us to conclude:

$$\begin{aligned}
 L(\mathbf{a^*b^*}) &= L(\mathbf{a^*})L(\mathbf{b^*}) \\
 &= \{uv \mid u \in L(\mathbf{a^*}) \wedge v \in L(\mathbf{b^*})\} \\
 &= \{uv \mid u \in \{\mathbf{a}^m \mid m \in \mathbb{N}\} \wedge v \in \{\mathbf{b}^n \mid n \in \mathbb{N}\}\} \\
 &= \{\mathbf{a}^m \mathbf{b}^n \mid m, n \in \mathbb{N}\}
 \end{aligned}$$

That is,  $L(\mathbf{a^*b^*})$  is the set of all words that start with a (possibly empty) sequence of  $\mathbf{a}$ 's, followed by a (possibly empty) sequence of  $\mathbf{b}$ 's.

- $(\epsilon + \mathbf{b})(\mathbf{ab})^*(\epsilon + \mathbf{a})$ :

Let us analyse the parts:

$$\begin{aligned}
 L(\epsilon + \mathbf{b}) &= \{\epsilon, \mathbf{b}\} \\
 L((\mathbf{ab})^*) &= \{(\mathbf{ab})^n \mid n \in \mathbb{N}\} \\
 L(\epsilon + \mathbf{a}) &= \{\epsilon, \mathbf{a}\}
 \end{aligned}$$

Thus, we have:

$$L((\epsilon + \mathbf{b})(\mathbf{ab})^*(\epsilon + \mathbf{a})) = \{u(\mathbf{ab})^n v \mid u \in \{\epsilon, \mathbf{b}\} \wedge n \in \mathbb{N} \wedge v \in \{\epsilon, \mathbf{a}\}\}$$

In English:  $L((\epsilon + \mathbf{b})(\mathbf{ab})^*(\epsilon + \mathbf{a}))$  is the set of (possibly empty) sequences of alternating  $\mathbf{a}$ 's and  $\mathbf{b}$ 's.



### 4.3 Algebraic laws

The semantics of regular expressions not only allows us to find out the meaning of specific regular expressions, but also allows us to prove useful laws about regular expression in general. Let us illustrate by proving the following distributive law for regular expressions:

$$E(F + G) = EF + EG$$

Note that  $E, F, G$  are *variables* standing for some specific but arbitrary regular expressions, and that  $=$  here is *semantic* (as opposed to syntactic) equality. That is, what we need to prove is that a regular expression of the form  $E(F + G)$  and one of the form  $EF + EG$  always have the same meaning, i.e., denote the same language.

We thus start from  $L(E(F + G))$  and show step by step that this is equal to  $L(EF + EG)$  without making any assumptions about the constituent regular expressions  $E, F$ , and  $G$ , other than that their semantics is given by  $L(E)$  etc.

$$\begin{aligned}
& L(E(F + G)) \\
= & \quad \{ \text{Semantics of r.e. (concat.)} \} \\
& L(E)L(F + G) \\
= & \quad \{ \text{Semantics of r.e. (+)} \} \\
& L(E)(L(F) \cup L(G)) \\
= & \quad \{ \text{Def. concat. of languages} \} \\
& \{w_1w_2 \mid w_1 \in L(E) \wedge w_2 \in (L(F) \cup L(G))\} \\
= & \quad \{ \text{Def. set union} \} \\
& \{w_1w_2 \mid w_1 \in L(E) \wedge w_2 \in \{w \mid w \in L(F) \vee w \in L(G)\}\} \\
= & \quad \{ \text{Duality between sets and predicates} \} \\
& \{w_1w_2 \mid w_1 \in L(E) \wedge (w_2 \in L(F) \vee w_2 \in L(G))\} \\
= & \quad \{ \text{Conjunction } (\wedge) \text{ distributes over disjunction } (\vee) \} \\
& \{w_1w_2 \mid (w_1 \in L(E) \wedge w_2 \in L(F)) \vee (w_1 \in L(E) \wedge w_2 \in L(G))\} \\
= & \quad \{ \text{Def. set union} \} \\
& \{w_1w_2 \mid (w_1 \in L(E) \wedge w_2 \in L(F))\} \cup \{w_1w_2 \mid (w_1 \in L(E) \wedge w_2 \in L(G))\} \\
= & \quad \{ \text{Def. concat. languages (twice)} \} \\
& L(E)L(F) \cup L(E)L(G) \\
= & \quad \{ \text{Semantics of r.e. (concat., twice)} \} \\
& L(EF) \cup L(EG) \\
= & \quad \{ \text{Semantics of r.e. (+)} \} \\
& L(EF + EG)
\end{aligned}$$

Other laws for regular expressions can be proved similarly, although induction is sometimes needed. As an exercise, prove (some of) the following:

$$\begin{aligned}
\epsilon E &= E \\
E\epsilon &= E \\
\emptyset E &= \emptyset \\
E\emptyset &= \emptyset \\
E + (F + G) &= (E + F) + G \\
E(FG) &= (EF)G \\
(E^*)^* &= E^* \\
\epsilon + EE^* &= E^*
\end{aligned}$$

## 4.4 Translating regular expressions into NFAs

**Theorem 4.1** *A regular expression  $E$  can be translated into an equivalent NFA  $N(E)$  such that  $L(N(E)) = L(E)$ .*

We refer to this translation as the “Graphical Construction”. It is a variation of the standard way of translating regular expressions into NFAs known as Thompson’s construction<sup>8</sup>.

**Proof.** The proof is by induction on the syntax of regular expressions:

1.  $N(\emptyset)$ :



which will reject everything (there are no final states). Thus:

$$\begin{aligned} L(N(\emptyset)) &= \emptyset \\ &= L(\emptyset) \end{aligned}$$

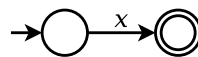
2.  $N(\epsilon)$ :



This automaton accepts the empty word but rejects everything else. Thus

$$\begin{aligned} L(N(\epsilon)) &= \{\epsilon\} \\ &= L(\epsilon) \end{aligned}$$

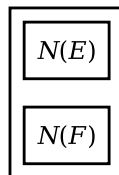
3.  $N(\mathbf{x})$ :



This automaton only accepts the word  $x$ . Thus:

$$\begin{aligned} L(N(\mathbf{x})) &= \{x\} \\ &= L(\mathbf{x}) \end{aligned}$$

4.  $N(E + F)$ : We merge the diagrams for  $N(E)$  and  $N(F)$  into one:



<sup>8</sup>[https://en.wikipedia.org/wiki/Thompson%27s\\_construction](https://en.wikipedia.org/wiki/Thompson%27s_construction)

Intuitively, the NFAs for the subexpressions  $E$  and  $F$  are placed side by side. Thus if either of the NFA accepts a word, the combined NFA accepts this word. However, we have to ensure that the states of the constituent NFAs do not get confused with each other. We therefore have to use the *disjoint union*, defined as follows:

$$A \uplus B = \{(0, a) \mid a \in A\} \cup \{(1, b) \mid b \in B\}$$

That is, each element of the disjoint union is “tagged” with an index that shows from which of the two sets it originated. Thus the elements of the constituent sets will remain distinct. The transition function of the combined NFA also has to be defined to work on tagged states.

Thus, given the NFAs for the subexpressions  $E$  and  $F$ :

$$\begin{aligned} N(E) &= (Q_E, \Sigma, \delta_E, S_E, F_E) \\ N(F) &= (Q_F, \Sigma, \delta_F, S_F, F_F) \end{aligned}$$

we construct the combined NFA for the regular expression  $E + F$ :

$$N(E + F) = (Q_{E+F}, \Sigma, \delta_{E+F}, S_{E+F}, F_{E+F})$$

where

$$\begin{aligned} Q_{E+F} &= Q_E \uplus Q_F \\ \delta_{E+F}((0, q), x) &= \{(0, q') \mid q' \in \delta_E(q, x)\} \\ \delta_{E+F}((1, q), x) &= \{(1, q') \mid q' \in \delta_F(q, x)\} \\ S_{E+F} &= S_E \uplus S_F \\ F_{E+F} &= F_E \uplus F_F \end{aligned}$$

It remains to prove  $L(N(E + F)) = L(E + F)$ . We first observe that

$$L(N(E + F)) = L(N(E)) \cup L(N(F))$$

because the initial states  $S_{E+F}$  of the combined NFA is the (disjoint) union of the initial states of the constituent NFAs, and because the combined NFA accepts a word whenever one of the constituent NFAs does.

The proof then proceeds by induction; that is, we assume that the translation is correct for the subexpressions:

$$\begin{aligned} L(N(E)) &= L(E) \\ L(N(F)) &= L(F) \end{aligned}$$

Thus:

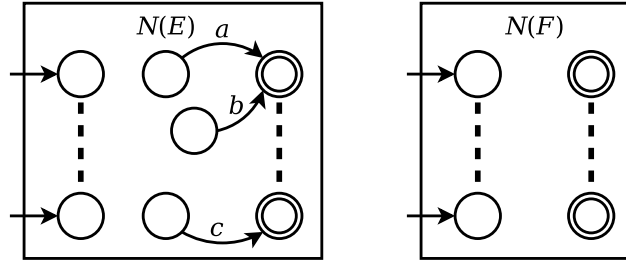
$$\begin{aligned} L(N(E + F)) &= L(N(E)) \cup L(N(F)) && \text{Above} \\ &= L(E) \cup L(F) && \text{Induction hypothesis} \\ &= L(E + F) && \text{By (4)} \end{aligned}$$

This is what is meant by *induction over the syntax of regular expressions*.

5.  $N(EF)$ : Recall that  $L(EF) = L(E)L(F)$ . Thus, a word  $w \in L(EF)$  iff  $w$  can be divided into two words  $u$  and  $v$ ,  $w = uv$ , such that  $u \in L(E)$  and  $v \in L(F)$ . Consequently, if we have an NFA  $N(E)$  recognising  $L(E)$  and another NFA  $N(F)$  recognising  $L(F)$ , we can construct an NFA recognising words  $w = uv \in L(EF)$  if we join the NFAs  $N(E)$  and  $N(F)$  in sequence in such a way that the machine moves from  $N(E)$  to an initial state of  $N(F)$  on the last symbol of a word  $u \in L(E)$ .

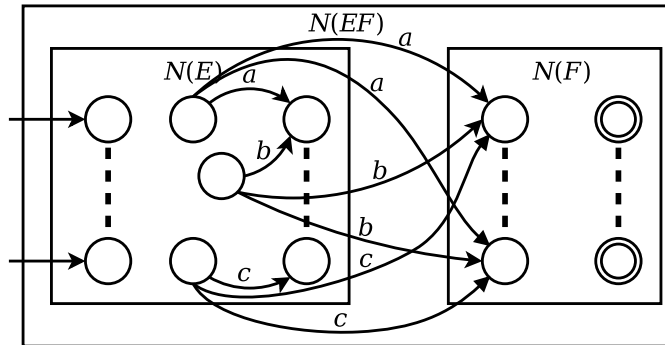
Of course, it could be that  $\epsilon \in L(E)$ , meaning that one or more of the initial states of  $N(E)$  are accepting. In this case, for a word  $w = uv$  with  $u = \epsilon$ , the machine needs to start in an initial state of  $N(F)$  directly as there is no last symbol of  $u = \epsilon$  to move on.

Therefore, we consider two cases:  $\epsilon \notin L(E)$ , meaning no initial state of  $N(E)$  is accepting, and  $\epsilon \in L(E)$ , meaning at least one initial state of  $N(E)$  is accepting. We start with the former:



The dashed lines here suggest “one or more”. So there could be one or more initial states and one or more final states in both  $N(E)$  and  $N(F)$ . This will be made precise shortly; the figure just conveys the general idea.

We identify every state of  $N(E)$  that *immediately precedes* an accepting state; i.e., every state from which an accepting state can be reached on a single input symbol. We then join  $N(E)$  and  $N(F)$  into a combined NFA  $N(EF)$  by adding an edge from each of the identified states to *all* of the initial states of  $N(F)$  for each symbol on which an accepting state of  $N(E)$  can be reached. The initial states of  $N(EF)$  are the initial states of  $N(E)$  and the final states of  $N(EF)$  are the final states of  $N(F)$ :



This ensures that the NFA  $N(EF)$ , once it has read a word  $u$  accepted by  $N(E)$ , is ready to try to accept the remainder  $v$  of a word  $w = uv$

by effectively passing  $v$  to  $N(F)$ , allowing the latter to try to accept the remaining part  $v$  of  $w$  from any of its initial states.

We now formalise this construction. The set of states of the combined NFA  $N(EF)$  is again given by the *disjoint* union of the states of  $N(E)$  and  $N(F)$  to avoid confusion, and the transition function  $\delta_{EF}$  as well as the initial states  $S_{EF}$  and final states  $F_{EF}$  are defined accordingly.

Thus, given the NFAs for the subexpressions  $E$  and  $F$ :

$$N(E) = (Q_E, \Sigma, \delta_E, S_E, F_E)$$

$$N(F) = (Q_F, \Sigma, \delta_F, S_F, F_F)$$

we construct the combined NFA for the regular expression  $EF$ :

$$N(EF) = (Q_{EF}, \Sigma, \delta_{EF}, S_{EF}, F_{EF})$$

where

$$Q_{EF} = Q_E \uplus Q_F$$

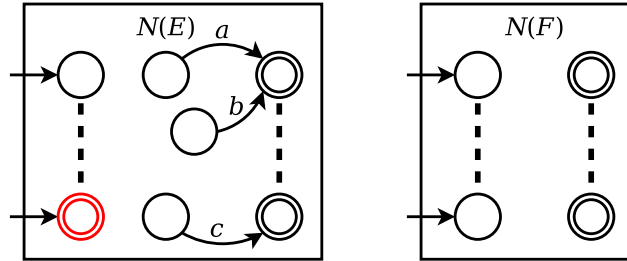
$$\begin{aligned} \delta_{EF}((0, q), x) &= \{(0, q') \mid q' \in \delta_E(q, x)\} \\ &\quad \cup \{(1, q') \mid \delta_E(q, x) \cap F_E \neq \emptyset \wedge q' \in S_F\} \end{aligned}$$

$$\delta_{EF}((1, q), x) = \{(1, q') \mid q' \in \delta_F(q, x)\}$$

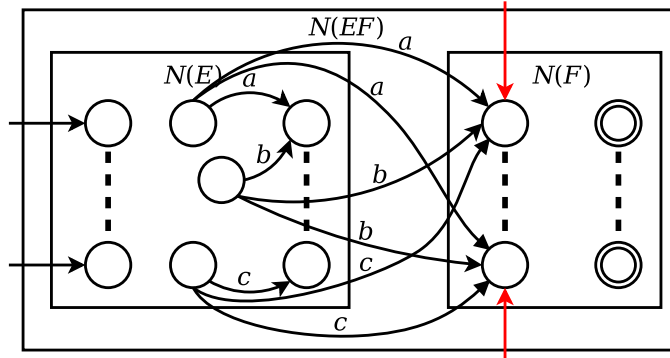
$$S_{EF} = \{(0, q) \mid q \in S_E\}$$

$$F_{EF} = \{(1, q) \mid q \in F_F\}$$

Now let us consider the second case: at least one of the initial states of  $N(E)$  is accepting:



As was discussed above, this simply means that we have to arrange that the initial states of  $N(F)$  also be initial states of the combined NFA  $N(EF)$ :



We thus refine the formal definition of the initial states of  $N(EF)$  to account for this, yielding a definition that covers both cases:

$$S_{EF} = \{(0, q) \mid q \in S_E\} \cup \{(1, q) \mid S_E \cap F_E \neq \emptyset \wedge q \in S_F\}$$

It remains to prove  $L(N(EF)) = L(EF)$ . From the construction above, it is clear that

$$L(N(EF)) = \{uv \mid u \in L(N(E)) \wedge v \in L(N(F))\}$$

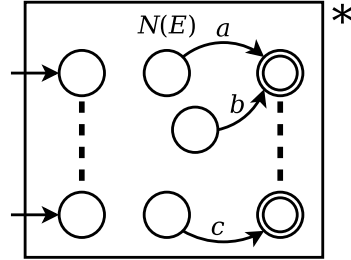
The proof again proceeds by induction; that is, we assume that the translation is correct for the subexpressions:

$$\begin{aligned} L(N(E)) &= L(E) \\ L(N(F)) &= L(F) \end{aligned}$$

Thus:

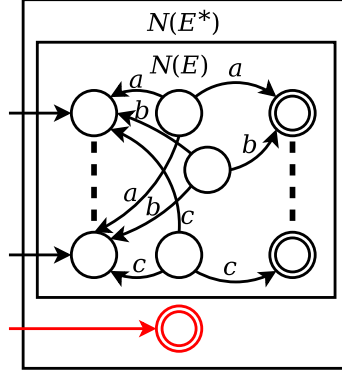
$$\begin{aligned} L(N(EF)) &= \{uv \mid u \in L(N(E)) \wedge v \in L(N(F))\} && \text{Above} \\ &= \{uv \mid u \in L(E) \wedge v \in L(F)\} && \text{Ind. hyp.} \\ &= L(E)L(F) && \text{Lang. concat.} \\ &= L(EF) && \text{By (5)} \end{aligned}$$

6.  $N(E^*)$ : Recall that  $L(E^*) = L(E)^*$ . Thus, a word  $w \in L(E^*)$  iff  $w$  can be divided into a sequence of  $n \in \mathbb{N}$  words  $u_i$ ,  $w = u_1u_2 \dots u_n$ , such that  $\forall i \in [1, n] . u_i \in L(E)$ . Consequently, if we have an NFA  $N(E)$  recognising  $L(E)$ , we can construct an NFA recognising words  $w \in L(E^*)$  by connecting 0 or more NFAs  $N(E)$  in sequence in a similar way to what we did for the case  $N(EF)$  above:



Here we use the  $*$  to informally suggest sequential composition of 0 or more NFAs.

However, we need to construct a *single* NFA, and there is no upper bound on the number of NFAs  $N(E)$  that we need to connect in sequence. We resolve this by taking a single NFA  $N(E)$  and construct an NFA for  $N(E^*)$  by making it loop back to all of its own initial states from each state that immediately precedes an accepting state. As we also need to allow for iteration 0 times, we further have to add one extra state that is both initial and final thus accepting  $\epsilon$ :



We now formalise this construction. This time, the states of  $N(E^*)$  are almost the same as those of  $N(E)$ . But we need one extra state to ensure that  $N(E^*)$  can accept the empty word,  $\epsilon$ , and we have to make sure that this one extra state cannot be confused with any other state in  $N(E)$ . We label the new state  $\epsilon$ , suggestive of its role to ensure acceptance of the empty word, and we then form the states of  $N(E^*)$  using disjoint union to ensure states cannot be accidentally confused. Like before, we have to take this into account when defining the transition function  $\delta_{E^*}$  as well as the initial states  $S_{E^*}$  and final states  $F_{E^*}$  of the NFA  $N(E^*)$ .

Thus, given the NFA resulting from translating the subexpression  $E$ :

$$N(E) = (Q_E, \Sigma, \delta_E, S_E, F_E)$$

we construct the NFA for the regular expression  $E^*$ :

$$N(E^*) = (Q_{E^*}, \Sigma, \delta_{E^*}, S_{E^*}, F_{E^*})$$

where

$$\begin{aligned} Q_{E^*} &= Q_E \uplus \{\epsilon\} \\ \delta_{E^*}((0, q), x) &= \{(0, q') \mid q' \in \delta_E(q, x)\} \\ &\quad \cup \{(0, q') \mid \delta_E(q, x) \cap F_E \neq \emptyset \wedge q' \in S_E\} \\ \delta_{E^*}((1, \epsilon), x) &= \emptyset \\ S_{E^*} &= S_E \uplus \{\epsilon\} \\ F_{E^*} &= F_E \uplus \{\epsilon\} \end{aligned}$$

It remains to prove  $L(N(E^*)) = L(E^*)$ . Given the construction above, we claim that

$$L(N(E^*)) = \{u_1 u_2 \dots u_n \mid n \in \mathbb{N} \wedge \forall i \in [1, n]. u_i \in L(N(E))\}$$

The intuition is that we can run through the automaton one or more times and that the new state  $\epsilon$  allows the NFA to accept the empty word.

The proof then again proceeds by induction; that is, we assume that the translation is correct for the subexpression:

$$L(N(E)) = L(E)$$

Thus:

$$\begin{aligned}
L(N(E^*)) &= \{ u_1 u_2 \dots u_n \mid n \in \mathbb{N} \wedge \forall i \in [1, n]. u_i \in L(N(E)) \} && \text{Above} \\
&= \{ u_1 u_2 \dots u_n \mid n \in \mathbb{N} \wedge \forall i \in [1, n]. u_i \in L(E) \} && \text{Ind. hyp.} \\
&= \bigcup_{n=0}^{\infty} L(E)^n && \text{Lang. concat.} \\
&= L(E)^* && \text{Def. Kleene star} \\
&= L(E^*) && \text{By (6)}
\end{aligned}$$

7.  $N((E)) = N(E)$ : Parentheses are just used for grouping and does not change anything.

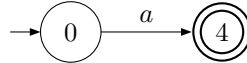
We need to prove  $L(N((E))) = L((E))$ . The proof is again by induction, so we assume  $L(N(E)) = L(E)$  and then we proceed as follows:

$$\begin{aligned}
L(N((E))) &= L(N(E)) && \text{By construction} \\
&= L(E) && \text{Induction hypothesis} \\
&= L((E)) && \text{By (7)}
\end{aligned}$$

□

It is worth pausing briefly to reflect on what we just have accomplished. In effect, we have implemented a *compiler* that translates regular expressions into NFAs, and we have proved it correct; that is, the translation preserves the meaning (here, the described language), which after all is what we generally expect of an accurate translation. Of course, it is a very simple compiler. Yet, in essence, it reflects how real tools that handle regular expressions work; for example, scanner (or lexer) generators such as Flex, Ragel, or Alex<sup>9</sup>. Moreover, while proving the correctness of compilers for typical programming languages is vastly more complicated than what we have seen here, there are methodological similarities, such as proof by induction over the structure of the language.

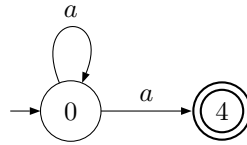
Let us illustrate how to apply the Graphical Construction. As a first example, we construct  $N(\mathbf{a^*b^*})$ . We start with the innermost subexpressions and then join the NFAs together step by step. The states are named according to how they will be named in the final NFA to make it easier to follow the derivation. It is fine to leave states unnamed until the end, and that is what normally is done. We begin with the NFA for  $\mathbf{a}$ :



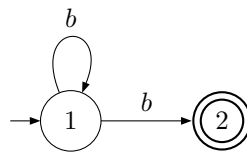
The NFA for  $\mathbf{a^*}$  is obtained by adding a loop on  $a$  from state 0 to itself as this state precedes a final state and is the only initial state, and by adding the extra state for accepting  $\epsilon$ :

<sup>9</sup>[https://en.wikipedia.org/wiki/Lexical\\_analysis](https://en.wikipedia.org/wiki/Lexical_analysis)

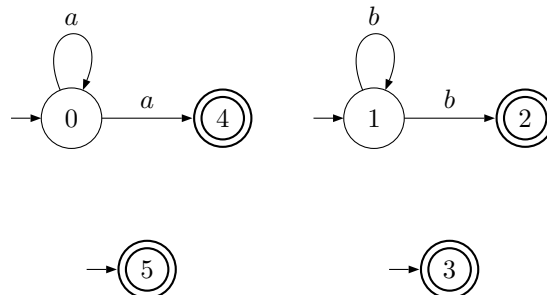




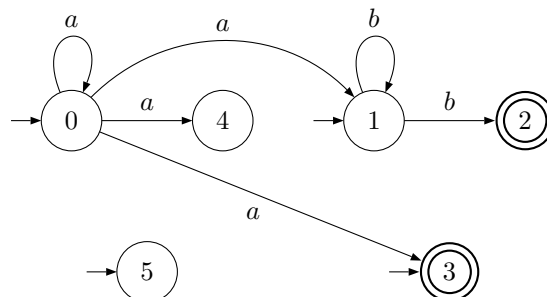
The NFA for  $\mathbf{b}^*$  is constructed in the same way:



Now we have to join these two NFAs in sequence:

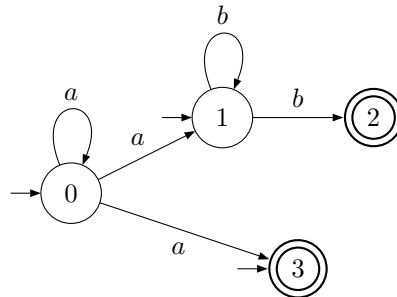


We have to pay extra attention because the automaton for the subexpression  $\mathbf{a}^*$  contains a state that is both initial and final, namely state 5, resulting in “extra” initial states when composing that automaton with the automaton for the subexpression  $\mathbf{b}^*$ :

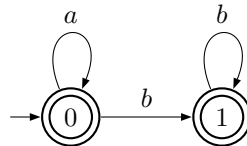


The states 4 and 5 have manifestly become “dead ends”: there is no way to reach a final state from either. For NFAs, such dead ends can simply be removed

(along with associated edges) without changing the accepted language. If we do that, we obtain:

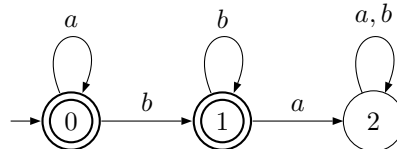


You may have noted that, though correct, this NFA is unnecessary complicated. For example, the following NFA also accepts  $N(\mathbf{a}^*\mathbf{b}^*)$ , but has fewer states:



This is typical: the translation of regular expressions into NFAs does generally not yield the simplest possible automata.

If we are interested in obtaining the smallest possible machine, one approach is to first convert the resulting NFA into a DFA using the subset construction of section 3.2.3 and then *minimize* this DFA as explained in section 5. If we do this for the four-state NFA above, we obtain the following DFA:

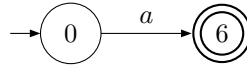


As it happens, just applying the subset construction to the four-state NFA yields this DFA directly<sup>10</sup>: it is already minimal. Try it! It is quick and a good exercise.

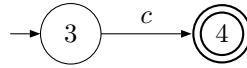
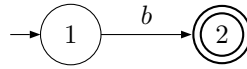
Let us do one more, somewhat larger, example: constructing an NFA for  $((\mathbf{a} + \epsilon)(\mathbf{b} + \mathbf{c}))^*$ . We again start with the innermost subexpressions and then join the NFAs together step by step. We have to (again) pay extra attention because the automaton for the subexpression  $(\mathbf{a} + \epsilon)$  contains a state that is both initial and final, resulting in “extra” initial states when composing that automaton with the automaton for the subexpression  $(\mathbf{b} + \mathbf{c})$ . Also, it makes sense to eliminate dead ends as soon as they occur, here *before* closing the loop due to the top-level Kleene star. The states are named according to how they will be named in the final NFA to make it easier to follow the derivation, but could be left unnamed until the end if you prefer.

First, an NFA for  $\mathbf{a} + \epsilon$ :

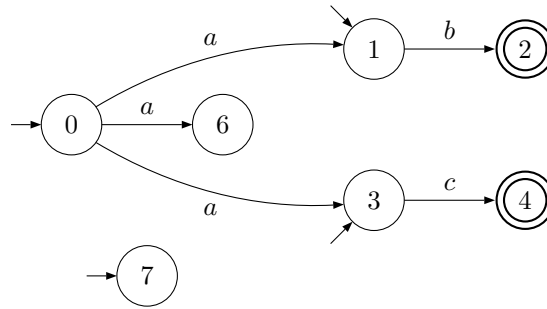
<sup>10</sup>Or one isomorphic to it: the states will probably be named differently.



NFA for  $\mathbf{b + c}$

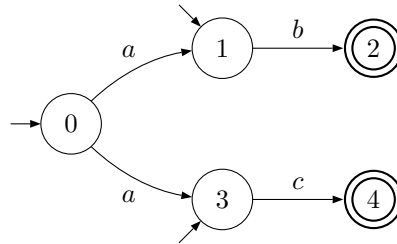


Join the above two NFAs to obtain an NFA for  $(\mathbf{a + \epsilon})(\mathbf{b + c})$ :

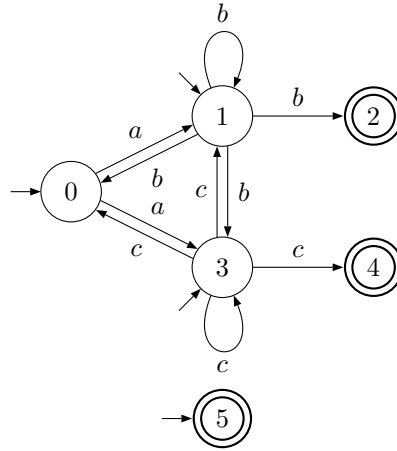


Note that both state 1 and 3 *remain* initial states because the left automaton has an initial state that is also accepting, meaning we need to be able to get to the start states of the right automaton without consuming any input.

States 6 and 7 have now manifestly become dead ends because there is no way to reach an accepting state from either. Let us remove them and all associated edges:



The last step is to carry out the construction corresponding to the  $*$ -operator. States 1 and 3 both immediately precede a final state, and we should thus add corresponding transition edges from those back to all initial states. There are three initial states, 0, 1, and 3. Thus we need an edge labelled  $b$  from 1 to each of 0, 1, and 3 (i.e., a loop back to itself on  $b$ ) and an edge labelled  $c$  from 3 to each of 0, 1, and 3 (i.e., a loop back to itself on  $c$ ). Additionally, we must not forget to add an extra initial state which is also final (here state 5) to ensure the NFA accepts  $\epsilon$ .



Note that the isolated state 5 thus also is part of the final automaton.

## 4.5 Summing up

From the previous section we know that a language given by regular expression is also recognized by a NFA. What about the other way: Can a language recognized by a finite automaton (DFA or NFA) also be described by a regular expression? The answer is yes:

**Theorem 4.2** *Given a DFA  $A$  there is a regular expression  $R(A)$  that recognizes the same language  $L(A) = L(R(A))$ .*

We omit the proof (which can be found in the [HMU01] on pp.91-93). However, we conclude:

**Corollary 4.3** *Given a language  $L \subseteq \Sigma^*$  the following is equivalent:*

1.  $L$  is given by a regular expression.
2.  $L$  is the language accepted by an NFA.
3.  $L$  is the language accepted by a DFA.

**Proof.** We have that 1.  $\implies$  2 by theorem 4.1. We know that 2.  $\implies$  3. by 3.2 and 3.  $\implies$  1. by 4.2.  $\square$

As indicated in the introduction, the languages that are characterised by any of the three equivalent conditions are called *regular languages* or *type 3 languages*.

## 4.6 Exercises

### Exercise 4.1

Give regular expressions defining the following languages over the alphabet  $\Sigma = \{a, b, c\}$ :

1. All words that contain exactly one  $a$ .
2. All words that contain at least two  $b$ s.
3. All words that contain at most two  $c$ s.
4. All words such that all  $b$ 's appear before all  $c$ 's.
5. All words that contain exactly one  $b$  and one  $c$  (but any number of  $a$ 's).
6. All words such that the number of  $a$ 's plus the number of  $b$ 's is odd.
7. All words that contain the sequence  $abba$  at least once.

### Exercise 4.2

Using the formal definition of the meaning of regular expressions, compute the set denoted by the regular expression

$$(\mathbf{aa} + \epsilon \mathbf{b}^* \emptyset)(\mathbf{b} + \mathbf{c})$$

simplifying as far as possible. Provide a step-by-step account.

### Exercise 4.3

Construct an NFA for the regular expression  $(\mathbf{a(b+c)})^*$  using the “graphical construction” from the lecture notes. Provide a step-by-step account.

For NFAs it is possible to omit “dead ends”, i.e., states from which no final state possibly can be reached, without changing the language of the automaton. Do this as soon as dead ends emerges to reduce your work.

### Exercise 4.4

Systematically construct an NFA for the regular expression

$$(\mathbf{a(\emptyset^* + b)})^*(\mathbf{c + \epsilon + \emptyset})$$

by following the graphical construction from the lecture notes. Make sure it is clear how you undertake the construction by showing the major steps. Eliminate “dead ends” (states from which no final state can be reached) when they appear. The states in the final NFA should be named, but as long as it is clear what you are doing, you can leave the states of intermediate NFAs unnamed.

## 5 Minimization of Finite Automata

As we saw when translating regular expressions into NFAs, the resulting automaton is not necessarily the smallest possible one. Similarly, when employing the subset construction to translate an NFA into a DFA, the result is not always the smallest possible DFA. It is often desirable to make automata as small as possible. For example, if we wish to implement an automaton, the implementation will be more efficient the smaller the automaton is.

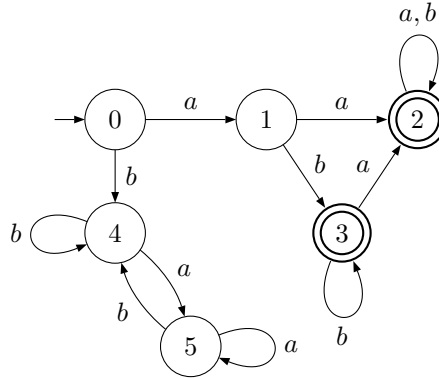
Given an automaton, the question, then, is how to construct an *equivalent* but smaller automaton. Recall that two automata are equivalent if they accept the same language. In the following we will study a method for minimizing DFAs: the table-filling algorithm.

Another interesting question is if there, in general, is one unique automaton that is the smallest equivalent one, or if there can be many distinct equivalent automata, none of which can be made any smaller. It turns out the answer is that the minimal equivalent DFA is unique up to naming of the states. This, in turn, means that we have obtained a mechanical decision procedure for determining whether two regular languages are equal: simply convert their respective representation (be it a DFA, an NFA, or a regular expression) to DFAs and minimize them. Because the minimal DFAs are unique, the languages are equal if and only if the minimal DFAs are equal.

### 5.1 The table-filling algorithm

For a DFA  $(Q, \Sigma, \delta, q_0, F)$ ,  $p, q \in Q$  are *equivalent* states if and only if, for all  $w \in \Sigma^*$ ,  $\hat{\delta}(p, w) \in F \Leftrightarrow \hat{\delta}(q, w) \in F$ . If two states are *not* equivalent, then they are *distinguishable*.

Consider the following DFA, where  $\Sigma = \{a, b\}$ ,  $Q = \{0, 1, 2, 3, 4, 5\}$ ,  $F = \{2, 3\}$ :



The states 1 and 2 are distinguishable on  $\epsilon$  because  $\hat{\delta}(1, \epsilon) = 1 \notin F$  while  $\hat{\delta}(2, \epsilon) = 2 \in F$ . Similarly, 0 and 1 are distinguishable on e.g.  $b$  because  $\hat{\delta}(0, b) = 4 \notin F$  while  $\hat{\delta}(1, b) = 3 \in F$ . On the other hand, in this case, we can easily see that 4 and 5 are not distinguishable on any word because it is not possible to reach any accepting (final) state from either 4 or 5.

The *Table-Filling Algorithm* recursively constructs the set of distinguishable pairs of states for a DFA. When all distinguishable state pairs have been

identified, any remaining pairs of states must be equivalent. Such states can be merged, thereby minimizing the automaton. Assume a DFA  $(Q, \Sigma, \delta, q_0, F)$ :

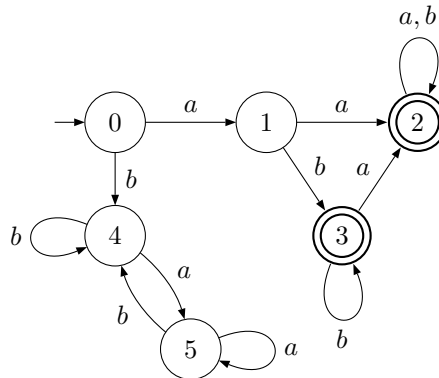
**BASIS** For  $p, q \in Q$ , if  $(p \in F \wedge q \notin F) \vee (p \notin F \wedge q \in F)$ , then  $(p, q)$  is a distinguishable pair of states. (The states  $p$  and  $q$  are distinguishable on  $\epsilon$ .)

**INDUCTION** For  $p, q, r, s \in Q$ ,  $a \in \Sigma$ , if  $(r, s) = (\delta(p, a), \delta(q, a))$  is a distinguishable pair of states, then  $(p, q)$  is also distinguishable. (If the states  $r$  and  $s$  are distinguishable on a word  $w$ , then  $p$  and  $q$  are distinguishable on  $aw$ .)

**Theorem 5.1** *If two states are not distinguishable by the table-filling algorithm, then they are equivalent.*

## 5.2 Example of DFA minimization using the table-filling algorithm

This section illustrates how the table-filling algorithm can be used to minimize a DFA when working by hand through a fully worked example. We will minimize the following DFA, where  $\Sigma = \{a, b\}$ ,  $Q = \{0, 1, 2, 3, 4, 5\}$ ,  $F = \{2, 3\}$ :



First construct a table over all pairs of distinct states. That is, we do not consider pairs  $(p \in Q, p \in Q)$  because a state obviously cannot be distinguishable from itself. An easy way of doing constructing the table is to order the states (e.g. numerically or alphabetically), and then list all states except the last one in *ascending* order along the top, and all states except the first one in *descending* order down the left-hand side of the table. The resulting table for our 6-state DFA looks like this:

	0	1	2	3	4
5					
4					
3					
2					
1					

Then mark the state pairs that are distinguishable according to the basis of the table-filling algorithm; i.e., the pairs where one state is accepting, and

one is not. The accepting states of our DFA are 2 and 3. Thus the state pairs  $(0, 2)$ ,  $(0, 3)$ ,  $(1, 2)$ ,  $(1, 3)$ ,  $(2, 4)$ ,  $(2, 5)$ ,  $(3, 4)$ , and  $(3, 5)$  have to be marked. List all remaining state pairs to the right: these are the potentially equivalent states that we now have to investigate further:

	0	1	2	3	4		<u>(0, 1)</u>	<u>(0, 4)</u>	<u>(0, 5)</u>	<u>(1, 4)</u>
5			x	x						
4			x	x						
3	x	x					<u>(1, 5)</u>	<u>(2, 3)</u>	<u>(4, 5)</u>	
2	x	x								
1										

Recall that the induction step of the table-filling algorithm says that for  $p, q, r, s \in Q$  and  $a \in \Sigma$ , if  $(r, s) = (\delta(p, a), \delta(q, a))$  is distinguishable (on some word  $w$ ), then  $(p, q)$  is (on the word  $aw$ ). If, during systematic investigation of all state combinations, we, from a state pair  $(p, q)$  on some input symbol  $a$ , reach a state pair  $(r, s)$  for which it is not yet known whether it is distinguishable or not, we record  $(p, q)$  under the heading for  $(r, s)$ . If it later becomes clear that  $(r, s)$  is distinguishable, that means that  $(p, q)$  also is distinguishable, and recording this implication allows us to carry out the deferred marking at that point.

Investigate all potentially equivalent state pairs on all input symbols (unless we find that a pair is distinguishable, which means we can stop):

- $(0, 1)$ :  $(\delta(0, a), \delta(1, a)) = (1, 2)$  Distinguishable! Mark in table.  
 $(0, 4)$ :  $(\delta(0, a), \delta(4, a)) = (1, 5)$  Unknown as yet. Add  $(0, 1)$  under  $(1, 5)$ .  
 $(\delta(0, b), \delta(4, b)) = (4, 4)$  Same state, no info.

Our table now looks as follows (we strike a line across the pairs we have considered):

	0	1	2	3	4		<del>-(0, 1)-</del>	<del>-(0, 4)-</del>	<u>(0, 5)</u>	<u>(1, 4)</u>
5			x	x						
4			x	x						
3	x	x					<u>(1, 5)</u>	<u>(2, 3)</u>	<u>(4, 5)</u>	
2	x	x					<u>(0, 4)</u>			
1	x									

We continue:

- $(0, 5)$ :  $(\delta(0, a), \delta(5, a)) = (1, 5)$  Unknown as yet. Add  $(0, 5)$  under  $(1, 5)$ .  
 $(\delta(0, b), \delta(5, b)) = (4, 4)$  Same state, no info.  
 $(1, 4)$ :  $(\delta(1, a), \delta(4, a)) = (2, 5)$  Distinguishable! Mark in table.

Table:

	0	1	2	3	4		<del>-(0, 1)-</del>	<del>-(0, 4)-</del>	<del>-(0, 5)-</del>	<del>-(1, 4)-</del>
5			x	x						
4		x	x	x						
3	x	x					<u>(1, 5)</u>	<u>(2, 3)</u>	<u>(4, 5)</u>	
2	x	x					<u>(0, 4)</u>			
1	x						<u>(0, 5)</u>			



Now we have come to the state pair  $(1, 5)$ . If we can determine that  $(1, 5)$  is a distinguishable pair, then we also know that the pairs  $(0, 4)$  and  $(0, 5)$  are distinguishable:

$(1, 5)$ :  $(\delta(1, a), \delta(5, a)) = (2, 5)$  Distinguishable!

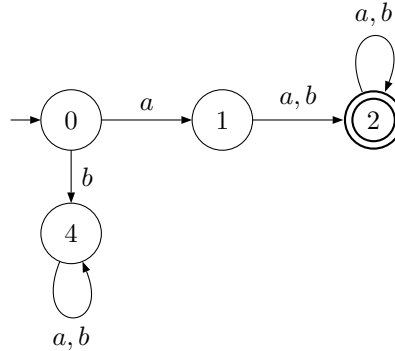
Thus we should mark  $(1, 5)$  along with  $(0, 4)$  and  $(0, 5)$ :

	0	1	2	3	4	
5	x	x	x	x		<del>-(0, 1)</del> <del>-(0, 4)</del> <del>-(0, 5)</del> <del>-(1, 4)</del>
4	x	x	x	x		
3	x	x				<del>-(1, 5)</del> <u>(2, 3)</u> <u>(4, 5)</u>
2	x	x				<del>-(0, 4)</del>
1	x					<del>-(0, 5)</del>

It remains to check the pairs  $(2, 3)$  and  $(4, 5)$ :

$(2, 3)$ :  $(\delta(2, a), \delta(3, a)) = (2, 2)$  Same state, no info.  
 $(\delta(2, b), \delta(3, b)) = (2, 3)$  No point in adding  $(2, 3)$  below  $(2, 3)$ .  
 $(4, 5)$ :  $(\delta(4, a), \delta(5, a)) = (5, 5)$  Same state, no info.  
 $(\delta(4, b), \delta(5, b)) = (4, 4)$  Same state, no info.

We have now systematically checked all potentially equivalent state pairs. Two pairs remain unmarked; i.e., we have not been able to show that they are distinguishable:  $(2, 3)$  and  $(4, 5)$ . We can therefore conclude that these states are pairwise equivalent:  $2 \equiv 3$  and  $4 \equiv 5$ . We thus proceed to merge these states by, informally, placing them “on top” of each other and “dragging along” the edges. The result is the following minimal DFA (where the merged states have been given the names 2 and 4):



## 6 Disproving Regularity

Regular languages are languages that can be recognized by a computer with finite memory. Such a computer corresponds to a DFA. However, there are many languages that cannot be recognized using only finite memory. A simple example is the language

$$L = \{0^n 1^n \mid n \in \mathbb{N}\}$$

That is, the language of words that start with a number of 0s followed by the *same* number of 1s. Note that this is different from  $L(0^*1^*)$ , which is the language of words of a sequences of 0s followed by a sequence of 1s but not necessarily of the same length. (We know this to be regular because it is given by a regular expression.)

Why can  $L$  not be recognized by a computer with finite memory? Suppose we have 32 MiB of memory; that is, we have  $32 * 1024 * 1024 * 8 = 268435456$  bits. Such a computer corresponds to an enormous DFA with  $2^{268435456}$  states (imagine drawing the transition diagram!). However, this computer can only count to  $2^{268435456} - 1$ ; if it reads a word starting with more than  $2^{268435456} - 1$  0s, it will necessarily lose count! The same reasoning applies whatever finite amount of memory we equip our computer with. Thus, an unbounded amount of memory is needed recognize  $L$ . (Of course,  $2^{268435456} - 1$  is a *very* large number indeed, so for practical purposes the machine will almost certainly be able to count much further than we ever will need.)

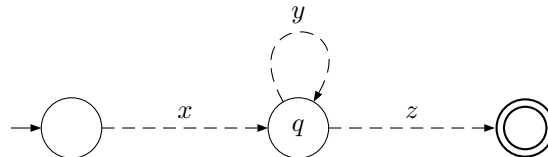
We will now show a general theorem called *the pumping lemma* (for regular languages) that allows us to prove that a certain language is not regular.

### 6.1 The pumping lemma

**Theorem 6.1** *Given a regular language  $L$ , then there is a number  $n \in \mathbb{N}$  such that all words  $w \in L$  that are longer than  $n$  ( $|w| \geq n$ ) can be split into three words  $w = xyz$  s.t.*

1.  $y \neq \epsilon$
2.  $|xy| \leq n$
3.  $\forall k \in \mathbb{N} . xy^kz \in L$

**Proof.** For a regular language  $L$  there exists a DFA  $A$  s.t.  $L = L(A)$ . Let us assume that  $A$  has  $n$  states. If  $A$  accepts a word  $w$  with  $|w| \geq n$ , it must have visited some state  $q$  twice:



We choose  $q$  such that it is the first cycle; hence  $|xy| \leq n$ . We also know that  $y$  is nonempty (otherwise there is no cycle). Now, consider what happens if a word of the form  $xy^kz$  is given to the automaton. The automaton will clearly accept this word; it just has to go round the cycle  $k$  times, whatever  $k$  is, including  $k = 0$ . Thus  $\forall k \in \mathbb{N} . xy^kz \in L$ .

□

## 6.2 Applying the pumping lemma

**Theorem 6.2** *The language  $L = \{0^n 1^n \mid n \in \mathbb{N}\}$  is not regular.*

**Proof.** Assume  $L$  would be regular. We will show that this leads to contradiction using the pumping lemma.

By the pumping lemma, there is an  $n$  such that we can split each word that is longer than  $n$  such that the properties given by the pumping lemma hold. Consider  $0^n 1^n \in L$ . This is certainly longer than  $n$ . We have that  $xyz = 0^n 1^n$  and we know that  $|xy| \leq n$ , hence  $y$  can only contain 0s. Further, because  $y \neq \epsilon$ , it must contain at least one 0. Now, according to the pumping lemma,  $xy^0 z \in L$ . However, this cannot be the case because it contains at least one fewer 0s than 1s. Our assumption that  $L$  is regular must thus have been wrong.  $\square$

It is easy to see that the language

$$\{1^n \mid n \text{ is even}\}$$

is regular (just construct the appropriate DFA or use a regular expression). However what about

$$\{1^n \mid n \text{ is a square}\}$$

where by saying  $n$  is a square we mean that there is an  $k \in \mathbb{N}$  s.t.  $n = k^2$ . We may try as we like: there is no way to find out whether we have a got a square number of 1s by only using finite memory. And indeed:

**Theorem 6.3** *The language  $L = \{1^n \mid n \text{ is a square}\}$  is not regular.*

**Proof.** We apply the same strategy as above. Assume  $L$  is regular. Then there is a number  $n$  such we can split all longer words according to the pumping lemma. Let us take  $w = 1^{n^2}$ ; this is certainly long enough. By the pumping lemma, we know that we can split  $w = xyz$  s.t. the conditions of the pumping lemma hold. In particular we know that

$$1 \leq |y| \leq |xy| \leq n$$

Using the 3rd condition we know that

$$xyyz \in L$$

that is  $|xyyz|$  is a square. However we know that

$$\begin{aligned} n^2 &= |w| \\ &= |xyz| \\ &< |xyyz| && \text{because } 1 \leq |y| \\ &= |xyz| + |y| \\ &\leq n^2 + n && \text{because } |y| \leq n \\ &< n^2 + 2n + 1 \\ &= (n+1)^2 \end{aligned}$$

To summarize, we have

$$n^2 < |xyyz| < (n+1)^2$$

That is  $|xyyz|$  lies between two subsequent squares. But then it cannot be a square itself, and hence we have a contradiction to  $xyyz \in L$ . We conclude  $L$  is not regular.  $\square$

Given a word  $w \in \Sigma^*$  we write  $w^R$  for the word read backwards. E.g.  $\text{abc}^R = \text{bca}$ . Formally this can be defined as

$$\begin{aligned}\epsilon^R &= \epsilon \\ (xw)^R &= w^R x\end{aligned}$$

We use this to define the language of even length palindromes

$$L_{\text{pali}} = \{ww^R \mid w \in \Sigma^*\}$$

E.g. for  $\Sigma = \{a, b\}$  we have  $\text{abba} \in L_{\text{pali}}$ . Using the intuition that finite automata can only use finite memory, it should be clear that this language is not regular either: to check whether the 2nd half is the same as the 1st half read backwards, we have to remember the first half, however long it is. Indeed, we can show:

**Theorem 6.4** *Given  $\Sigma = \{a, b\}$  we have that  $L_{\text{pali}}$  is not regular.*

**Proof.** We use the pumping lemma: We assume that  $L_{\text{pali}}$  is regular. Now given a pumping number  $n$  we construct  $w = \text{a}^n \text{bba}^n \in L_{\text{pali}}$ , this word is certainly longer than  $n$ . From the pumping lemma we know that there is a splitting of the word  $w = xyz$  s.t.  $|xy| \leq n$  and hence  $y$  may only contain  $a$ 's and because  $y \neq \epsilon$  at least one. We conclude that  $xz \in L_{\text{pali}}$  where  $xz = \text{a}^m \text{bba}^n$  where  $m < n$ . However, this word is not a palindrome, because the sequence of  $a$ 's at the beginning is shorter than the sequence of  $a$ 's at the end. Hence our assumption  $L_{\text{pali}}$  is regular must be wrong.  $\square$

The proof works for any alphabet with at least 2 different symbols. However, if  $\Sigma$  contains only one symbol, as in  $\Sigma = \{1\}$ , then  $L_{\text{pali}}$  is the language of an even number of 1s and this is regular:  $L_{\text{pali}} = (11)^*$ .

## 6.3 Exercises

### Exercise 6.1

Apply the pumping lemma for regular languages to show that the following languages are not regular:

1.  $L_1 = \{a^n b^m c^{n+m} \mid m, n \in \mathbb{N}\}$  over the alphabet  $\Sigma_1 = \{a, b, c\}$ .  
(E.g.,  $\text{aabbcccc} \in L_1$ , but  $\text{aabbcc} \notin L_1$ .)
2.  $L_2 = \{w \in \Sigma^* \mid \#_a(w) = 2 \times \#_b(w) \wedge \#_b(w) = 2 \times \#_c(w)\}$  over the alphabet  $\Sigma = \{a, b, c\}$ , where  $\#_a(w)$  denotes the number of  $a$ 's in a word  $w$ ,  $\#_b(w)$  the number of  $b$ 's, etc.  
(E.g.,  $\text{abcaaba} \in L_2$ , but  $\text{aaabbc} \notin L_2$ , and  $\text{aaaaabbc} \notin L_2$ .)

## 7 Context-Free Grammars

This section introduces *context-free grammars* (CFGs) as a formalism to define languages that is more general than regular expressions; that is, there are more languages definable by CFGs than by regular expressions and finite automata. The class of languages definable by CFGs is known as the *context-free languages* or *type 2 languages* (section 1.1). We will define the notion of automata corresponding to CFGs, the *push down automata* (PDA), later, in section 9.

Context-free grammars have an abundance of applications. A prominent example is the definition of (aspects of) the syntax programming languages, like C, Java, or Haskell. Another application is the document type definition (DTD) for the SGML-family markup languages, like XML and HTML. Applications of CFGs are discussed further in section 7.6.

### 7.1 What are context-free grammars?

We start by defining what context-free grammars are, their *syntax*, deferring what CFGs mean, the languages they describe or their *semantics*, to section 7.2. A context-free grammar  $G = (N, T, P, S)$  is given by

- A finite set  $N$  of *nonterminal symbols* or *nonterminals*.
- A finite set  $T$  of *terminal symbols* or *terminals*.
- $N \cap T = \emptyset$ ; i.e., the sets  $N$  and  $T$  are disjoint.
- A finite set  $P \subseteq N \times (N \cup T)^*$  of productions. A production  $(A, \alpha)$ , where  $A \in N$  and  $\alpha \in (N \cup T)^*$  is a sequence of nonterminal and terminal symbols. It is written as  $A \rightarrow \alpha$  in the following.
- $S \in N$ : the distinguished start symbol.

Nonterminals are also referred to as *variables* and consequently the set of nonterminals is sometimes denoted by  $V$ . Yet another term for the same thing is *syntactic categories*. The terminals are the alphabet of the language defined by a CFG, and for that reason the set of terminals is sometimes denoted by  $\Sigma$ . Indeed, we will occasionally use that convention as well in the following.

Note that the right-hand side of a production may be empty. This is known as an  $\epsilon$ -production and written

$$A \rightarrow \epsilon$$

Also note that it is perfectly permissible for the *same* nonterminal to occur both to the left and to the right of the arrow in a production. In fact, this is essential: if that were not allowed, CFGs would only amount to a (possibly) compact description of finite languages and be of little interest. A production for a nonterminal  $A$  where the same nonterminal is the first symbol of the right-hand side, in the leftmost position, is called *immediately left-recursive*; e.g.,

$$A \rightarrow A\alpha$$

where  $\alpha \in (N \cup T)^*$ . A production for a nonterminal  $A$  where the same nonterminal is the last symbol of the right-hand side, in the rightmost position, is called *immediately right-recursive*; e.g.,

$$A \rightarrow \alpha A$$

Recursion can also be *indirect*: the left-hand side non-terminal of a production can be reached again from the right-hand side via one or more other productions. This is very common: see the following example for an illustration.

As an example we define a grammar for the language of arithmetic expressions over  $a$  using only  $+$  and  $*$ . As we will see in section 7.2 where the language described by a CFG is defined, the elements of this language are words like  $a + (a * a)$  or  $(a + a) * (a + a)$ . On the other hand, words like  $a ++a$  or  $)(a$ , which manifestly do not correspond to well-formed arithmetic expressions, do not belong to the language:

$$G_{\text{arith}} = (\{E, T, F\}, \{ (, ), a, +, * \}, P, E)$$

where  $P$  is given by:

$$\begin{aligned} P = \{ & E \rightarrow T, \\ & E \rightarrow E + T, \\ & T \rightarrow F, \\ & T \rightarrow T * F, \\ & F \rightarrow a, \\ & F \rightarrow (E) \} \end{aligned}$$

Here, the choice  $E, T, F$  for the nonterminal symbols is meant to suggest *Expression*, *Term*, and *Factor*, respectively. Note that some of the productions for  $E$  and  $T$  are immediately left-recursive, and that one of the productions for  $F$  recursively refer back to the start symbol  $S$ . The latter is an example of indirect recursion.

Somewhat unfortunately,  $T$  is used here both as one of the nonterminals and to denote the set of terminals, as per the conventions outlined above. Do not let this confuse you: the nonterminal symbol  $T$  and the set of terminal symbols are quite distinct! Occasional name clashes are a fact of life.

To save space, we may combine all the rules with the same left-hand side, separating the alternatives with a vertical bar. Using this convention, our set of productions can be written

$$\begin{aligned} P = \{ & E \rightarrow T \mid E + T, \\ & T \rightarrow F \mid T * F, \\ & F \rightarrow a \mid (E) \} \end{aligned}$$

In practice, the set of productions is often given by just listing the productions, without explicit braces indicating a set:

$$\begin{aligned} E &\rightarrow T \mid E + T \\ T &\rightarrow F \mid T * F \\ F &\rightarrow a \mid (E) \end{aligned}$$

Either way, these are just a more convenient ways to write down exactly the same set of productions.

## 7.2 The meaning of context-free grammars

How can we check if a word  $w \in T^*$  is in the language of a grammar? We start with the start symbol  $S$ . Note that this is a word in  $(N \cup T)^*$ . If there is a production  $S \rightarrow \alpha$ , we can obtain a new word in  $(N \cup T)^*$  by replacing the  $S$  by the right-hand side  $\alpha$  of the production. Should there be further nonterminals in the resulting word, the process is repeated by looking for a production where the left-hand side is one of those nonterminals and replacing that nonterminal by the right-hand side of the production. This process is called a *derivation*. It is often the case that there is a choice between derivation steps as there can be more than one nonterminal in a word and more than one production for a nonterminal. Any word  $w \in T^*$  derived in this way belongs to the language defined by the grammar.

Let us consider our expression grammar  $G_{\text{arith}}$  from section 7.1. In this case, the start symbol is  $E$ . The following is one possible derivation:

$$\begin{aligned}
 E &\Rightarrow E + T \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow T + T \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow F + T \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + T \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + F \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (E) \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (T) \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (T * F) \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (F * F) \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (a * F) \\
 &\quad G_{\text{arith}} \\
 &\Rightarrow a + (a * a) \\
 &\quad G_{\text{arith}}
 \end{aligned}$$

Generally, given a grammar  $G$ , the symbol  $\Rightarrow_G$  stands for the relation *derives in one step in grammar  $G$*  or *directly derives in grammar  $G$* . It has nothing to do with implication. When the grammar used is clear from the context, it is conventional to drop the subscript  $G$  and simply use  $\Rightarrow$ , read *directly derives*. In the example above, we always replaced the leftmost nonterminal symbol. This is called a *leftmost* derivation. However, as was remarked before, this is not necessary: in general, we are free to pick any nonterminal for replacement. An alternative would be to always pick the rightmost one, which results in a *rightmost* derivation. Or we could pick whichever nonterminal seems more convenient to expand. The symbols  $\Rightarrow_{\text{lm}}$  and  $\Rightarrow_{\text{rm}}$  are sometimes used to indicate leftmost and rightmost derivation steps, respectively.

Given any grammar  $G = (N, T, P, S)$  we define the relation *directly derives in grammar  $G$*  as follows:

$$\begin{aligned}
 \Rightarrow_G &\subseteq (N \cup T)^* \times (N \cup T)^* \\
 \alpha A \gamma \Rightarrow_G \alpha \beta \gamma &\iff A \rightarrow \beta \in P
 \end{aligned}$$

The relation *derives in grammar*  $G$ , derivation in 0 or more steps, is defined as:

$$\begin{aligned} \xRightarrow{*}_G &\subseteq (N \cup T)^* \times (N \cup T)^* \\ \alpha_0 \xRightarrow{*}_G \alpha_n &\iff \alpha_0 \xRightarrow{G} \alpha_1 \xRightarrow{G} \dots \alpha_{n-1} \xRightarrow{G} \alpha_n \end{aligned}$$

where  $n \in \mathbb{N}$ . Thus  $\alpha \xRightarrow{*}_G \alpha$ , as  $n$  may be 0. We will also occasionally use  $\xRightarrow{+}_G$  meaning derivation in one or more steps<sup>12</sup>.

A word  $\alpha \in (N \cup T)^*$  such that  $S \xRightarrow{*}_G \alpha$  is called a *sentential form*. The *language* of a grammar,  $L(G) \subseteq T^*$ , consists of all terminal sentential forms:

$$L(G) = \{w \in T^* \mid S \xRightarrow{*}_G w\}$$

A language that can be defined by a context-free grammar is called a context-free language (CFL).

### 7.3 The relation between regular and context-free languages

A grammar in which each production has at most one non-terminal symbol in its right-hand side is *linear*. For example,

$$G_1 = (\{S\}, \{0, 1\}, \{S \rightarrow \epsilon \mid 0S1\}, S)$$

is a linear grammar. There are two special cases of linear grammars:

- A linear grammar is *left-linear* if each right-hand side nonterminal is the leftmost (first) symbol in its right-hand side.
- A linear grammar is *right-linear* if each right-hand side nonterminal is the rightmost (last) symbol in its right-hand side.

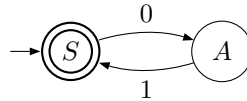
For example,

$$G_2 = (\{S, A\}, \{0, 1\}, \{S \rightarrow \epsilon \mid A1, A \rightarrow S0\}, S)$$

is left-linear, and

$$G_3 = (\{S, A\}, \{0, 1\}, \{S \rightarrow \epsilon \mid 0A, A \rightarrow 1S\}, S)$$

is right-linear. Collectively, left-linear and right-linear grammars are called *regular* grammars because the languages they describe are *regular*. We will not prove this fact, but it is easy to see how right-linear grammars correspond directly to NFAs. For example,  $G_3$  corresponds to the following NFA:



<sup>12</sup>  $\xRightarrow{*}_G$  is the *reflexive transitive closure* of  $\xRightarrow{G}$ , and  $\xRightarrow{+}_G$  is the *transitive closure* of  $\xRightarrow{G}$ .



Thus we see that context-free grammars can be used to describe at least some regular languages. On the other hand, some of the languages that we have shown not to be regular are actually context-free. For example, the (linear) grammar  $G_1$  above describes the language  $\{0^n 1^n \mid n \in \mathbb{N}\}$  that we proved (theorem 6.2) not to be regular. It is worth noting that if we allow left-linear and right-linear productions to be mixed, the resulting language is not necessarily regular. For example, the grammar

$$G'_1 = (\{S, A\}, \{0, 1\}, \{S \rightarrow \epsilon \mid 0A, A \rightarrow S1\}, S)$$

is equivalent to  $G_1$ ; i.e., describes the same, non-regular, language. Further, we proved (theorem 6.4) that the language of even-length palindromes

$$L_{\text{pali}} = \{ww^R \mid w \in \{a, b\}^*\}$$

is not regular. The following context-free grammar is one way of defining this language, demonstrating that  $L_{\text{pali}}$  is a context-free language:

$$G_{\text{pali}} = (\{S\}, \{a, b\}, \{S \rightarrow \epsilon \mid aSa \mid bSb\}, S)$$

So, what is the relation between the regular and context-free languages? Are there languages that are regular but not context-free? The answer is no:

**Theorem 7.1** *All regular languages are context-free.*

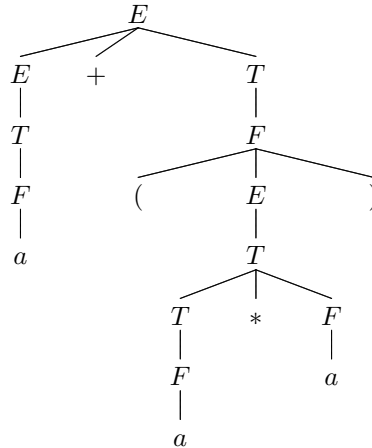
Again, we do not give a proof, but the idea is that regular expressions can be translated into (regular) context-free grammars. For example,  $\mathbf{a^*b^*}$  can be translated into:

$$(\{A, B\}, \{a, b\}, \{A \rightarrow aA \mid B, B \rightarrow bB \mid \epsilon\}, A)$$

As a consequence of theorem 7.1 and the fact that we have seen that there are at least some languages that are context-free but not regular, we have established that the regular languages form a *proper subset* of the context-free ones.

## 7.4 Derivation trees

A derivation in a context-free grammar induce a corresponding *derivation tree* that reflects the *structure* of the derivation: how each nonterminal was rewritten. As an example, consider the tree representation of the derivation of  $a + (a * a)$  in grammar  $G_{\text{arith}}$  (section 7.2):



The central point is that the parent-child relationship reflects a derivation step according to a production in the grammar. For example, the production  $F \rightarrow (E)$  was used once in the derivation of  $a + (a * a)$ , and consequently there is a node labelled  $F$  in the derivation tree with child nodes labelled  $(, E, )$ , ordered from left to right.

In more detail, a tree is a *derivation tree* for a CFG  $G = (N, T, P, S)$  iff:

1. Every node has a label from  $N \cup T \cup \{\epsilon\}$ .
2. The label of the root node is  $S$ .
3. Labels of internal nodes belong to  $N$ .
4. If a node  $n$  has label  $A$  and nodes  $n_1, n_2, \dots, n_k$  are children of  $n$ , from left to right, with labels  $X_1, X_2, \dots, X_k$ , respectively, then  $A \rightarrow X_1 X_2 \dots X_k$  is a production in  $P$ .
5. If a node  $n$  has label  $\epsilon$ , then  $n$  is a leaf and the only child of its parent.

Through the notion of the *yield* of a derivation tree, the relationship between a derivation tree and corresponding derivations can be made precise:

- The string of *leaf labels* read from left to right, eliding any  $\epsilon$  bar one if it is the only remaining symbol, constitute the *yield* of the tree.
- For a CFG  $G = (N, T, P, S)$ , a string  $\alpha \in (N \cup T)^*$  is the yield of some derivation tree iff  $S \xRightarrow[G]{*} \alpha$ .

Note that the leaf nodes may be labelled with either terminal or nonterminal symbols. The yield is thus a *sentential form* (see section 7.2) in general, and not necessarily a word in the language of the context-free grammar.

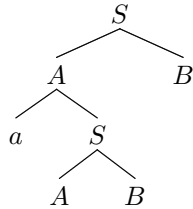
To illustrate the above point, along with elision of superfluous  $\epsilon$ s from the yield, consider the grammar

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aS \mid \epsilon \\ B &\rightarrow Sb \mid \epsilon \end{aligned}$$

One derivation in this grammar is

$$S \Rightarrow AB \Rightarrow aSB \Rightarrow aABB$$

The corresponding derivation tree is

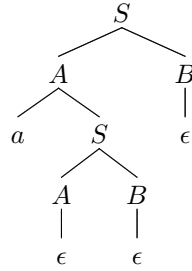


Here, the yield is the sentential form  $aABB$ .

We can continue the derivation, using the  $\epsilon$ -productions for the nonterminals  $A$  and  $B$ :

$$aABB \Rightarrow aBB \Rightarrow aB \Rightarrow a$$

The derivation tree for the entire derivation is

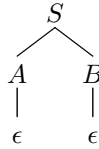


with the yield being just  $a$ .

For an example where the yield is empty, consider the derivation

$$S \Rightarrow AB \Rightarrow B \Rightarrow \epsilon$$

The derivation tree is



with yield  $\epsilon$ .

## 7.5 Ambiguity

A CFG  $G = (N, T, P, S)$  is *ambiguous* iff there is at least one word  $w \in L(G)$  such that there are

- two different *derivation trees*, or equivalently
- two different *leftmost derivations*, or equivalently
- two different *rightmost derivations*

for  $w$ . This is usually a bad thing because it entails that there is more than one way to interpret a word; i.e., it leads to semantic ambiguity.

As an example consider the following variation of a grammar for simple arithmetic expressions (SAE):

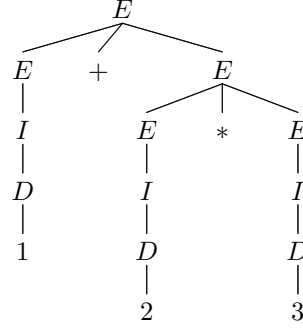
$$SAE = (N = \{E, I, D\}, T = \{+, *, (, ), 0, 1, \dots, 9\}, P, E)$$

where  $P$  is given by:

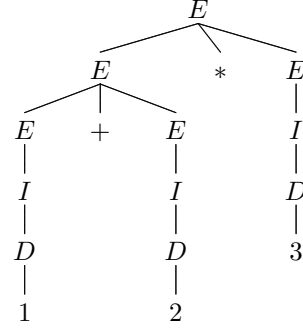
$$\begin{aligned}
 E &\rightarrow E + E \\
 &\quad | \quad E * E \\
 &\quad | \quad (E) \\
 &\quad | \quad I \\
 I &\rightarrow DI \mid D \\
 D &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9
 \end{aligned}$$

The grammar  $SAE$  allows expressions involving numbers to be derived, unlike  $G_{arith}$ . Also note that this grammar is simpler in that there only is one nonterminal ( $E$ ) involved at the level of expressions proper, as opposed to three ( $E$ ,  $T$ ,  $F$ ) for  $G_{arith}$ .

Consider the word  $1 + 2 * 3$ . The following derivation tree shows that this word belongs to  $L(SAE)$ :



Note that the yield is the word above,  $1 + 2 * 3$ . But there is *another* tree with the same yield:



Thus, there is *one* word for which there are *two* derivation trees. This shows that the grammar  $SAE$  is ambiguous.

As per the definition of ambiguity, another way to demonstrate ambiguity is to find two leftmost or two rightmost derivations for a word. It is easy to see that there is a one-to-one correspondence between a derivation tree and a leftmost and a rightmost derivation. To illustrate, here are the two leftmost derivations for  $1 + 2 * 3$ , corresponding to the first and to the second tree respectively:

$$\begin{aligned}
 E &\Rightarrow_{lm} E + E \Rightarrow_{lm} I + E \Rightarrow_{lm} D + E \Rightarrow_{lm} 1 + E \Rightarrow_{lm} 1 + E * E \\
 &\Rightarrow_{lm} 1 + I * E \Rightarrow_{lm} 1 + D * E \Rightarrow_{lm} 1 + 2 * E \\
 &\Rightarrow_{lm} 1 + 2 * I \Rightarrow_{lm} 1 + 2 * D \Rightarrow_{lm} 1 + 2 * 3
 \end{aligned}$$

and

$$\begin{aligned}
 E &\Rightarrow_{lm} E * E \Rightarrow_{lm} E + E * E \Rightarrow_{lm} I + E * E \Rightarrow_{lm} D + E * E \\
 &\Rightarrow_{lm} 1 + E * E \Rightarrow_{lm} 1 + I * E \Rightarrow_{lm} 1 + D * E \Rightarrow_{lm} 1 + 2 * E \\
 &\Rightarrow_{lm} 1 + 2 * I \Rightarrow_{lm} 1 + 2 * D \Rightarrow_{lm} 1 + 2 * 3
 \end{aligned}$$

Note, *one* word, *two* different *leftmost* derivations. Thus the grammar is ambiguous. Exercise: Find the *rightmost* derivation corresponding to each tree.

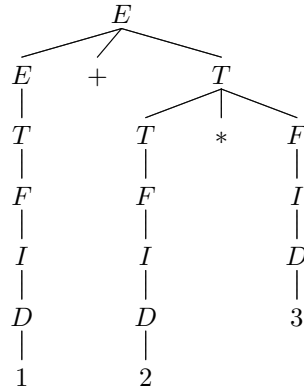
There are two reasons for why ambiguity is problematic. The first we already alluded to: semantic ambiguity. Suppose we wish to assign a meaning to a word beyond the word itself. In our case, the meaning might be the result of evaluating the expression, for instance. Then the first tree suggests that the expression should be read as  $1 + (2 * 3)$ , which evaluates to 7, while the second tree suggests the expression should be read as  $(1 + 2) * 3$ , which evaluates to 9<sup>13</sup>. Note how the bracketing reflects the tree structure in each case. We have *one* word with *two* different *interpretations* and thus semantic ambiguity.

The other reason is that many methods for *parsing*, i.e. determining if a given word belongs to the language defined by a CFG, do not work for ambiguous grammars. In particular, this applies to efficient methods for parsing that usually are what we would like to use for that very reason. We return to parsing in section 10.

Fortunately, it is often possible to change an ambiguous grammar into an unambiguous one that is equivalent; i.e., the language is unchanged. Such *grammar transformations* are discussed 8. But for now, let us note that  $G_{\text{arith}}$  was carefully structured so as to be unambiguous. We can restructure  $SAE$  similarly:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid I \\ I &\rightarrow DI \mid D \\ D &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

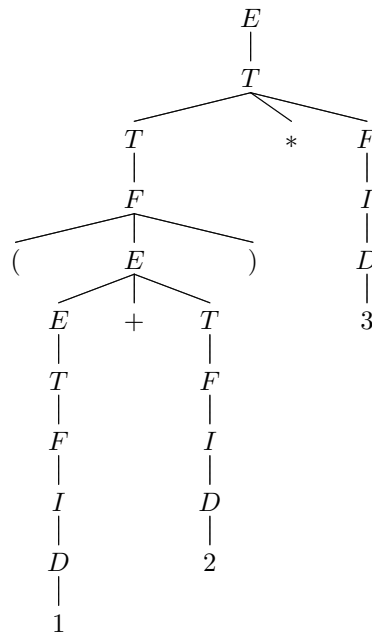
Now there is only one possible derivation tree for  $1 + 2 * 3$ :



Convince yourself that this is the case. Note that this tree corresponds to the reading  $1 + (2 * 3)$  of the expression; i.e. a reading where multiplication has higher precedence than addition. This is, of course, is the standard convention.

<sup>13</sup>Strictly speaking, this is just a natural and often convenient convention, because it implies that the meaning of the whole can be understood in terms of the meaning of the parts in the obvious way implied by the structure of the tree. However, sometimes, in real compilers, it might be necessary or convenient to parse things in a way that does not reflect the semantics as directly. But then the parse tree is usually transformed later into a simplified, internal version (an *Abstract Syntax Tree*), the structure of which reflects the intended semantics as described here.

Now, if the interpretation  $(1 + 2) * 3$  is desired, *explicit* parentheses have to be used. The (only!) derivation tree for this word is:



Note that the addition now is a subexpression (subtree) of the overall expression, which is what was desired.

## 7.6 Applications of context-free grammars

An important example for context-free languages is the syntax of programming languages. For example, the specification of the Java programming language<sup>14</sup> [GJS<sup>+</sup>15] uses a context-free grammar to specify the syntax of Java. Another example is the specification of the syntax of the language used in the G53CMP Compilers module [Nil16].

In practice, the exact details of how a grammar is presented often differs a little bit from the conventions introduced here. The Java language specification being a case in point. Further common variations include Backus-Naur form (BNF)<sup>15</sup> and Extended Backus-Naur form (EBNF)<sup>16</sup>. But these differences are entirely superficial. For example, the symbol  $::=$  is often used instead of  $\rightarrow$ . For another example, EBNF introduces a shorthand notation for “zero or more of”, commonly denoted by enclosing a grammar symbol (or symbols) in curly braces or using a Kleene-star like notation; e.g.:

$$A \rightarrow \dots \{B\} \dots$$

or

$$A \rightarrow \dots B^* \dots$$

<sup>14</sup><http://docs.oracle.com/javase/specs/>

<sup>15</sup>[https://en.wikipedia.org/wiki/Backus%E2%80%93Naur\\_form](https://en.wikipedia.org/wiki/Backus%E2%80%93Naur_form)

<sup>16</sup>[https://en.wikipedia.org/wiki/Extended\\_Backus%E2%80%93Naur\\_form](https://en.wikipedia.org/wiki/Extended_Backus%E2%80%93Naur_form)

It is easy to see that this iterative construct really just is a shorthand for basic productions as it readily can be expressed by introducing an auxiliary nonterminal and a couple of associated productions, either using left recursion or right recursion. The example above can be translated into the equivalent immediately left-recursive (section 7.1) productions

$$\begin{aligned} A &\rightarrow \dots A_1 \dots \\ A_1 &\rightarrow \epsilon \mid A_1 B \end{aligned}$$

or into the equivalent immediately right-recursive productions

$$\begin{aligned} A &\rightarrow \dots A_1 \dots \\ A_1 &\rightarrow \epsilon \mid B A_1 \end{aligned}$$

where  $A_1$  is a the new, auxiliary, nonterminal.

Note that not all syntactic aspects of common programming languages are captured by the context-free grammar. For example, requirements regarding declaring variables before they are used and type correctness of expressions cannot be captured through context-free languages. However, at least in the area of programming languages, it is common practice to use the notion of “syntactically correct” in the more limited sense of “conforming to the context-free grammar of the language”. Aspects such as type-correctness are considered separately.

Another application is the document type definition (DTD)<sup>17</sup> for the SGML-family markup languages, like XML and HTML. A DTD defines the document structure including legal elements and attributes. It can be declared inline, as a header of the document to which the definition applies, or as an external reference.

Extensions of context-free grammars are used in computer linguistics to describe natural languages. In fact, as mentioned in section 1.1, context-free grammars were originally invented by Noam Chomsky for describing natural languages. As an example, consider the following set of terminals, each an English word:

$$T = \{\text{the, dog, cat, that, bites, barks, catches}\}$$

We can then define a grammar  $G = (\{S, N, NP, VI, VT, VP\}, T, P, S)$  where  $P$  is the following set of productions:

$$\begin{aligned} S &\rightarrow NP VP \\ N &\rightarrow \text{cat} \mid \text{dog} \\ NP &\rightarrow \text{the } N \mid NP \text{ that } VP \\ VI &\rightarrow \text{barks} \mid \text{bites} \\ VT &\rightarrow \text{bites} \mid \text{catches} \\ VP &\rightarrow VI \mid VT NP \end{aligned}$$

This grammar allows us to derive interesting sentences like:

the dog that catches the cat that bites barks

---

<sup>17</sup>[https://en.wikipedia.org/wiki/Document\\_type\\_definition](https://en.wikipedia.org/wiki/Document_type_definition)

## 7.7 Exercises

### Exercise 7.1

Consider the following Context-Free Grammar (CFG)  $G$ :

$$\begin{aligned} S &\rightarrow X \mid Y \\ X &\rightarrow aXb \mid \epsilon \\ Y &\rightarrow cYd \mid \epsilon \end{aligned}$$

$S, X, Y$  are nonterminal symbols,  $S$  is the start symbol, and  $a, b, c, d$  are terminal symbols.

1. Derive the following words using the grammar  $G$ . Answer by giving the entire derivation sequence from the start symbol  $S$ :
  - (a)  $\epsilon$
  - (b)  $aabb$
  - (c)  $ccddd$
2. Does the string  $aaadd$  belong to the language  $L(G)$  generated by the grammar  $G$ ? Provide a brief justification.
3. Give a set expression (using set comprehensions and operations on sets like union) denoting the language  $L(G)$ .

### Exercise 7.2

Construct a context free grammar generating the language

$$L = \{(ab)^m(bc)^n(cb)^n(ba)^m \mid m, n \geq 1\} \cup \{d^n \mid n \geq 0\} \cup \{d^n \mid n \geq 2\}$$

over the alphabet  $\{a, b, c, d\}$  (parentheses are only used for grouping). Note that set concatenation has higher precedence than set union. Explain your construction.

### Exercise 7.3

Consider the Context-Free Grammar (CFG)  $G = (N, T, P, S)$  where  $N = \{S, X, Y\}$  are the nonterminal symbols,  $T = \{a, b, c\}$  are the terminal symbols,  $S$  is the start symbol, and the set of productions  $P$  is:

$$\begin{aligned} S &\rightarrow X \mid Y \\ X &\rightarrow aXb \mid ab \\ Y &\rightarrow bYc \mid \epsilon \end{aligned}$$

Recall that the relation “derives directly in  $G$ ” is a relation on strings of terminals and non-terminals; i.e.

$$\Rightarrow_G \subseteq (N \cup T)^* \times (N \cup T)^*$$

such that for all  $\alpha, \gamma \in (N \cup T)^*$

$$\alpha A \gamma \Rightarrow_G \alpha \beta \gamma \iff A \rightarrow \beta \in P$$

List all pairs  $(\phi, \theta)$  of the relation  $\Rightarrow_G$  for the cases where either  $\phi \in \{X, XY, aXbYc, cc\}$  or  $\theta = a$ .



### Exercise 7.4

Consider the following Context-Free Grammar (CFG) *Exp*:

$$\begin{aligned}
 T &\rightarrow T + T \mid F \\
 F &\rightarrow F * F \mid P \\
 P &\rightarrow N(A) \mid (T) \mid I \\
 N &\rightarrow f \mid g \mid h \\
 A &\rightarrow T \mid \epsilon \\
 I &\rightarrow DI \mid D \\
 D &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9
 \end{aligned}$$

$T, F, P, N, A, I, D$  are nonterminals;  $+, *, f, g, h, (, ), 0, 1, 2, 3, 4, 5, 6, 7, 8, 9$  are terminals;  $T$  is the start symbol.

- Derive the following words in the grammar *Exp* where possible. If it is possible to derive the word, give the entire *left-most* derivation; i.e. always expand the left-most non-terminal of the sentential form<sup>18</sup>. If it is not possible to derive the word, give a brief explanation as to why not.
  - (789)
  - $7 + g(3 * 5) * (f())$
  - $1 + 2 * 3$
  - $1 + 7(9)$
- Draw a derivation tree for the word  $7 + (8 * h(1)) + 9$  in the grammar *Exp*.
- Draw another derivation tree for the word  $7 + (8 * h(1)) + 9$  from 2. What does the fact that there are two different derivation trees for one word tell about the grammar *Exp*?
- Modify the relevant productions of the grammar *Exp* so that a function symbol (one of  $f, g, h$ ) can be applied to zero, one, or more arguments, separated by a single comma when there are more than one argument, instead of just zero or one argument. For example, it should be possible to derive words like

$$f(2, g(), h(3 + 4))$$

Explain your construction.

---

<sup>18</sup>*Sentential form*: word derivable from the start symbol.

## 8 Transformations of context-free grammars

In this section, we discuss how context free grammars can be restructured systematically without changing the language that they describe. There are many reasons for doing this, such as simplifying a grammar, putting it into some specific form, or eliminating ambiguities.

### 8.1 Equivalence of context-free grammars

Two grammars  $G_1$  and  $G_2$  are *equivalent* iff

$$L(G_1) = L(G_2)$$

Whenever a grammar is being transformed, it is assumed that the resulting grammar is equivalent to the original one.

For example, the following two grammars are equivalent:

$$\begin{array}{ll} G_1: & \begin{array}{l} S \rightarrow \epsilon \mid A \\ A \rightarrow a \mid aA \end{array} \end{array} \qquad \begin{array}{ll} G_2: & \begin{array}{l} S \rightarrow A \\ A \rightarrow \epsilon \mid Aa \end{array} \end{array}$$

$L(G_1) = \{a\}^* = L(G_2)$ . The equivalence of CFGs is in general *undecidable* (section 11).

### 8.2 Elimination of useless productions

In a context free grammar  $G = (N, T, P, S)$ , a nonterminal  $X \in N$  is

- *reachable* iff  $S \xRightarrow{*} \alpha X \beta$  for some  $\alpha, \beta \in \{N \cup T\}^*$
- *productive* iff  $X \xRightarrow{*} w$  for some  $w \in T^*$

Phrased differently, a nonterminal  $X$  is reachable if it occurs in some sentential form. Productions for nonterminals that are unreachable or unproductive, or where an unproductive nonterminal occurs on the right-hand side, are collectively known as *useless productions*. Any useless production can be removed from a grammar without changing the language.

For example, consider the following grammar, where  $N = \{S, A, B\}$ ,  $T = \{a, b\}$ , and  $S$  is the start symbol:

$$\begin{array}{l} S \rightarrow aAB \mid b \\ A \rightarrow aA \mid a \\ B \rightarrow bB \end{array}$$

The nonterminal  $B$  is unproductive as there is no way to derive a word of only terminal symbols from it. This makes the productions  $S \rightarrow aAB$  and  $B \rightarrow bB$  useless. Removing those productions leaves us with:

$$\begin{array}{l} S \rightarrow b \\ A \rightarrow aA \mid a \end{array}$$

But now  $A$  has clearly become an unreachable nonterminal, making the productions  $A \rightarrow aA$  and  $A \rightarrow a$  useless too. If we remove those as well we obtain:

$$S \rightarrow b$$

Of course, any terminal and nonterminals that no longer occur in any productions can also be eliminated. Thus the set of nonterminals is now just  $\{S\}$  and the set of terminals just  $\{b\}$ .

### 8.3 Substitution

As a direct consequence of how derivation in a grammar is defined, it follows that an occurrence of a non-terminal in a right-hand side may be replaced by the right-hand sides of the productions for that non-terminal if done in all possible ways. This is a form of *substitution*.

For example, consider the following grammar fragment. We assume it includes *all* productions for  $B$ , and we wish to eliminate the occurrence of  $B$  in the right-hand side of the production for  $A$ :

$$\begin{aligned} A &\rightarrow XBY \\ B &\rightarrow C \mid D \\ B &\rightarrow \epsilon \end{aligned}$$

Note that there are three productions for  $B$ . By substitution, this grammar fragment can be transformed into:

$$\begin{aligned} A &\rightarrow XCY \mid XDY \mid XY \\ B &\rightarrow C \mid D \\ B &\rightarrow \epsilon \end{aligned}$$

Thus we get one new production for  $A$  for each alternative for  $B$ . Note that we cannot necessarily remove the productions for  $B$ : there may be other occurrences of  $B$  in the grammar (as the above was only a fragment). However, if substitution renders some symbols unreachable, then the productions for those symbols become useless and can consequently be removed (section 8.2).

Substitution can of course also be performed on recursive productions. That is often not very useful, though, as the productions just get bigger and more numerous without eliminating any nonterminals. For example, assume the following productions are *all* productions for  $A$  in some grammar. Note that  $A \rightarrow Ab$  is recursive (immediately left-recursive, as it happens):

$$A \rightarrow \epsilon \mid Ab$$

We can substitute the right-hand sides of all productions for  $A$  for the one occurrence of  $A$  in the right-hand side of  $A \rightarrow Ab$ . That would leave us with:

$$A \rightarrow \epsilon \mid b \mid Abb$$

Thus we did not gain anything here, unless there was some specific reason for wanting a single production that shows that  $A$  can yield a single  $b$ .

## 8.4 Left factoring

Sometimes it is useful to factor out a common prefix of the right-hand sides of a group of productions. This is called *left factoring*. Left factoring can make a grammar easier to read and understand as it captures a recurring pattern in one place. It is sometimes also necessary for putting a grammar into a form suitable for use with certain parsing methods (section 10).

Consider the following two productions for  $A$ . Note the common prefix  $XY$ :

$$A \rightarrow XYX \mid XYZZY$$

After left factoring:

$$\begin{aligned} A &\rightarrow BX \mid BZZY \\ B &\rightarrow XY \end{aligned}$$

## 8.5 Disambiguating context-free grammars

Given an ambiguous context-free grammar  $G$ , it is often possible to construct an equivalent but *unambiguous* grammar  $G'$ . Some context-free languages are *inherently ambiguous*, meaning that every CFG generating the language is necessarily ambiguous, but most languages of practical interest, such as programming languages, can be given unambiguous CFGs.

In this section, we focus on how expression languages (like arithmetic expressions or regular expressions) can be given unambiguous CFGs by structuring the grammar to account for *operator precedence* and *operator associativity*.

The following is a CFG for simple arithmetic expressions. For simplicity, we only consider the numbers 0, 1, and 2:

$$\begin{aligned} E &\rightarrow E + E \mid E * E \mid E \uparrow E \mid (E) \mid N \\ N &\rightarrow 0 \mid 1 \mid 2 \end{aligned}$$

$E$  and  $N$  are nonterminals,  $E$  is the start symbol,  $+$ ,  $*$ ,  $\uparrow$ ,  $(, )$ ,  $0$ ,  $1$ ,  $2$  are terminals. The grammar is ambiguous. For example, the word  $0 + 1 + 2$  has two different leftmost derivations:

$$\begin{aligned} E &\Rightarrow_{\text{lm}} E + E \Rightarrow_{\text{lm}} N + E \Rightarrow_{\text{lm}} 0 + E \Rightarrow_{\text{lm}} 0 + E + E \Rightarrow_{\text{lm}} 0 + N + E \Rightarrow_{\text{lm}} 0 + 1 + E \\ &\Rightarrow_{\text{lm}} 0 + 1 + N \Rightarrow_{\text{lm}} 0 + 1 + 2 \\ E &\Rightarrow_{\text{lm}} E + E \Rightarrow_{\text{lm}} E + E + E \Rightarrow_{\text{lm}} N + E + E \Rightarrow_{\text{lm}} 0 + E + E \Rightarrow_{\text{lm}} 0 + N + E \\ &\Rightarrow_{\text{lm}} 0 + 1 + E \Rightarrow_{\text{lm}} 0 + 1 + N \Rightarrow_{\text{lm}} 0 + 1 + 2 \end{aligned}$$

We now wish to construct an equivalent but unambiguous version of the the above grammar by making it reflect the following conventions regarding operator precedence and associativity:

Operators	Precedence	Associativity
$\uparrow$	highest	right
$*$	medium	left
$+$	lowest	left

For example, the word

$$1 + 2 * 2 \uparrow 2 \uparrow 2 + 0$$

should be read as

$$(1 + (2 * (2 \uparrow (2 \uparrow 2)))) + 0$$

That is, the structure of the (one and only) derivation tree should reflect this reading.

To impart operator precedence on the grammar, it has to be *stratified*: First the expressions have to be partitioned into different categories of expressions, one category for each operator precedence level. We thus need to introduce one nonterminal (or syntactic category) for each precedence level. Then it must be arranged so that expressions belonging to the category for operators of one precedence level only occur as *subexpressions* of expressions belonging to the category for operators of the next lower precedence level. Bracketing (enclosing an expression in some form of parentheses) should have higher precedence than any operator. Bracketed subexpressions should thus only be allowed as subexpressions of expressions in the category for the highest operator precedence. The expression enclosed in the parentheses, however, should be of the category for the lowest operator precedence.

In our example, there are three operator precedence levels, so we introduce three additional nonterminals ( $E_1$ ,  $E_2$ ,  $E_3$ ) to stratify the grammar into four levels: one level for each operator precedence category, and one more for the innermost expression level, the subexpressions of expressions involving operators of the highest precedence. The resulting grammar is:

$$\begin{aligned} E &\rightarrow E_1 + E_1 \mid E_1 \\ E_1 &\rightarrow E_2 * E_2 \mid E_2 \\ E_2 &\rightarrow E_3 \uparrow E_3 \mid E_3 \\ E_3 &\rightarrow (E) \mid N \\ N &\rightarrow 0 \mid 1 \mid 2 \end{aligned}$$

This grammar is fine in that it is unambiguous. But for practical purposes, it is inconvenient as expressions involving more than one operator at the same level has to be explicitly bracketed. For example, the word  $0 + 1 + 2$  cannot be derived in this grammar (try it!), but a user of this little expression language would have to write either  $(0 + 1) + 2$  or  $0 + (1 + 2)$ .

This is where operator associativity comes into the picture: by adopting a convention regarding associativity, expressions involving more than operator of some specific precedence level can be implicitly bracketed. In our example,  $+$  is left-associative (as per standard mathematical conventions), which means the word  $0 + 1 + 2$  should be read as  $(0 + 1) + 2$ ; i.e. the derivation tree for  $0 + 1 + 2$  should branch to the left, suggesting that  $0 + 1$  is a subexpression of the overall expression.

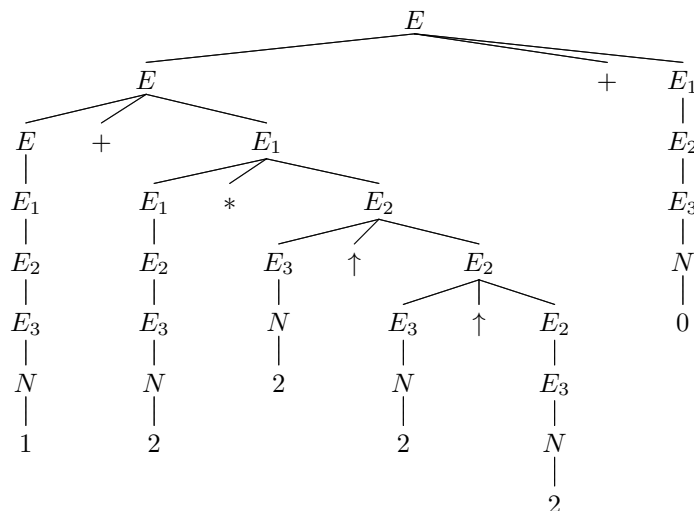
Imparting of associativity is achieved by making productions for *left*-associative operators *left*-recursive, and those for *right*-associative operators *right*-recursive, as this results in the desired left-branching or right-branching structure, respectively, of the derivation tree. Note that the requirement that subexpressions at one precedence level must all be expressions at next higher precedence level is relaxed a little, but only in a way that does not make the grammar unambiguous.

In our example, the operators  $+$  and  $*$  are left-associative, and the operator  $\uparrow$  is right-associative. We thus make the corresponding productions left- and right-recursive, respectively, arriving at the final version of the grammar:

$$\begin{array}{lcl} E & \rightarrow & E + E_1 \mid E_1 \\ E_1 & \rightarrow & E_1 * E_2 \mid E_2 \\ E_2 & \rightarrow & E_3 \uparrow E_2 \mid E_3 \\ E_3 & \rightarrow & (E) \mid N \\ N & \rightarrow & 0 \mid 1 \mid 2 \end{array}$$

This grammar is (again) unambiguous, but now allows words like  $0 + 1 + 2$ . Verify that there only is one derivation tree, and that this has the desired left-branching structure!

As an example, let us consider the word  $1 + 2 * 2 \uparrow 2 \uparrow 2 + 0$  from above. Recall that we want the reading  $(1 + (2 * (2 \uparrow (2 \uparrow 2)))) + 0$ . The derivation tree for the word in the final version of the grammar is:



Make sure you understand how the branching structure corresponds to the desired reading of the word; i.e., how precedence and associativity allowed for implicit bracketing. As an exercise, draw the derivation tree for the *explicitly* bracketed word, and compare the structure of the two trees. The trees will not be the same, of course, as the words are not the same (as words, symbol by symbol), but the branching structure (what is a subexpression of what) will agree.

## 8.6 Elimination of left recursion

A CFG is *left-recursive* if there is some non-terminal  $A$  such that  $A \xrightarrow{+} A\alpha$ <sup>19</sup>. Certain parsing methods cannot handle left-recursive grammars. An example is recursive decent parsing as described in section 10. If we want to use such a parsing method for parsing a language  $L = L(G)$  given by a left-recursive

<sup>19</sup>Recall that  $\stackrel{+}{\Rightarrow}$  means derives in one or more steps; i.e., the transitive closure of  $\Rightarrow$ .

grammar  $G$ , then it first has to be transformed into an equivalent grammar  $G'$  that is *not* left-recursive.

We first consider *immediate* left recursion (section 7.1); i.e., productions of the form  $A \rightarrow A\alpha$ . We assume that  $\alpha$  cannot derive  $\epsilon$ .

The key idea of the transformation is simple. Let us first consider a simplified scenario, with one immediately left-recursive production for a nonterminal  $A$  and one non-recursive production:

$$\begin{aligned} A &\rightarrow A\alpha \\ A &\rightarrow \beta \end{aligned}$$

Then observe that all strings derivable from  $A$  using these two productions have the form  $\beta(\alpha)^*$  once the last  $A$  has been replaced by  $\beta$ ; i.e.,  $\beta$  followed by zero or more  $\alpha$ . Now it is easy to see that the following alternative grammar generates strings of exactly the same form:

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \epsilon \end{aligned}$$

This is arguably a more direct way to generate strings of the form  $\beta(\alpha)^*$ . In essence, the productions say: “start with  $\beta$ , and then tag on zero or more  $\alpha$ ”.

We now generalise and formalise this idea. In order to transform an immediately left-recursive grammar to an equivalent grammar that is not left recursive, proceed as follows. For each nonterminal  $A$  defined by some left-recursive production, group the productions for  $A$

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

such that no  $\beta_i$  begins with an  $A$ . Then replace the  $A$  productions by

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{aligned}$$

Assumption: no  $\alpha_i$  derives  $\epsilon$ .

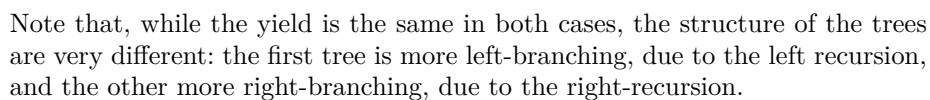
To illustrate, consider the immediately left-recursive grammar

$$\begin{aligned} S &\rightarrow A \mid B \\ A &\rightarrow ABc \mid AAdd \mid a \mid aa \\ B &\rightarrow Bee \mid b \end{aligned}$$

Applying the transformation rules above yields the following equivalent right-recursive grammar:

$$\begin{aligned} S &\rightarrow A \mid B & B &\rightarrow bB' \\ A &\rightarrow aA' \mid aaA' & B' &\rightarrow eeB' \mid \epsilon \\ A' &\rightarrow BcA' \mid AddA' \mid \epsilon \end{aligned}$$

Let us do a sanity check on the new grammar by picking a word in the language of the original grammar and make sure it can also be derived in the new grammar. The derivation of the word should ideally make use of all recursive productions in the grammar. Let us pick the word *aabeeecddbeec*. The following is the derivation tree for this word in the original grammar, demonstrating that the word indeed is in the language generated by the grammar:



To eliminate *general* left recursion, the grammar is first transformed into an *immediately* left-recursive grammar through systematic substitution (section 8.3). After this the elimination scheme set out above can be applied.

Consider the following grammar. Note that no production is immediately left-recursive:

The grammar is, however, left recursive because, for example,  $A \Rightarrow BaB \Rightarrow CbaB \Rightarrow AbbaB$ . Thus we have  $A \xRightarrow{+} A\alpha$  (for  $\alpha = bbaB$  in this case), demonstrating that the grammar is left recursive.



To eliminate the left recursion, let us first transform this grammar into an equivalent immediately left-recursive grammar. Let us start by eliminating  $C$  by substituting all alternatives for  $C$  into the right-hand side of the production  $B \rightarrow Cb$ . Note that this makes  $C$  an unreachable nonterminal, making the productions for  $C$  useless meaning they can be removed (section 8.2):

$$\begin{aligned} A &\rightarrow BaB \\ B &\rightarrow Abb \mid Acb \mid \epsilon \end{aligned}$$

Then we can eliminate  $B$  wherever it occurs in the *leftmost* position of a right-hand side in the productions for  $A$ . The grammar is now immediately left-recursive:

$$\begin{aligned} A &\rightarrow AbbaB \mid AcbaB \mid aB \\ B &\rightarrow Abb \mid Acb \mid \epsilon \end{aligned}$$

Alternatively,  $B$  can be eliminated completely from the productions for  $A$ , making  $B$  an unreachable terminal allowing the productions for  $B$  to be removed:

$$\begin{aligned} A &\rightarrow AbbaAbb \mid AcbaAbb \mid aAbb \\ &\mid AbbaAcb \mid AcbaAcb \mid aAcb \\ &\mid Abba \mid Acba \mid a \end{aligned}$$

Let us go with the smaller version (fewer productions):

$$\begin{aligned} A &\rightarrow AbbaB \mid AcbaB \mid aB \\ B &\rightarrow Abb \mid Acb \mid \epsilon \end{aligned}$$

Only the productions for  $A$  are immediately left-recursive. Applying the transformation to eliminate left recursion gives us:

$$\begin{aligned} A &\rightarrow aBA' \\ A' &\rightarrow bbaBA' \mid cbaBA' \mid \epsilon \\ B &\rightarrow Abb \mid Acb \mid \epsilon \end{aligned}$$

Note that  $A$  appears to the left in  $B$ -productions; yet the grammar is no longer left-recursive. Why?

## 8.7 Exercises

### Exercise 8.1

The grammar *Exp* from exercise 7.4 is ambiguous. Fix this problem; i.e., modify the grammar, *without* changing the language of the grammar, so that all words in the language have exactly one derivation tree. You do not need to prove that this holds for the resulting grammar, but you should explain what you did and why.

### Exercise 8.2

- Construct a simple, *unambiguous* grammar according to the following:
  - The integer literals are the only primitive expressions. An integer literal is either 0, or a non-empty word of decimal digits (0, 1, ..., 9) not starting with 0 and with a single optional minus sign (−) in front. E.g. 0, 1, 42, −234 are all valid integer literals, but 01, −0, − − 1 are not.
  - There are four binary operators:

Operators	Precedence	Associativity
<	1 (lowest)	non-associative
$\oplus$	2	left
$\otimes$	3	left
$\uparrow$	4 (highest)	right

- Additionally, it should be possible to use parentheses for grouping in the standard way.
- Draw the derivation tree for  $42 < 0 \otimes -10 \otimes (1 \oplus 7) \uparrow 2$  and verify that its structure reflects the specification above.

### Exercise 8.3

The following context-free grammar is immediately left-recursive:

$$\begin{aligned}
 S &\rightarrow Sa \mid XbS \mid a \\
 X &\rightarrow XXX \mid YYY \mid XYY \mid YYX \\
 Y &\rightarrow cY \mid dY \mid e
 \end{aligned}$$

$S$ ,  $X$ , and  $Y$  are nonterminals,  $S$  is the start symbol, and  $a$ ,  $b$ ,  $c$ ,  $d$ , and  $e$  are terminals.

Transform it into an equivalent *right-recursive* grammar. Explicitly show the result of the *grouping step* in addition to the final result after applying the actual *transformation step* to the productions that need transformation.

## 9 Pushdown Automata

We will now consider a new notion of automata *Pushdown Automata* (PDA). PDAs are finite automata with a stack; i.e., a data structure that can be used to store an arbitrary number of symbols (hence PDAs have an infinite set of states) but which can be only accessed in a *last-in-first-out* (LIFO) fashion. The languages that can be recognized by PDA are precisely the context-free languages.

### 9.1 What is a pushdown automaton?

A Pushdown Automaton  $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$  is given by the following data

- A finite set  $Q$  of states,
- A finite set  $\Sigma$  of input symbols (the alphabet),
- A finite set  $\Gamma$  of stack symbols,
- A transition function

$$\delta \in Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow \mathcal{P}_{\text{fin}}(Q \times \Gamma^*)$$

Here  $\mathcal{P}_{\text{fin}}(A)$  are the finite subsets of a set; i.e., this can be defined as

$$\mathcal{P}_{\text{fin}}(A) = \{X \mid X \subseteq A \wedge X \text{ is finite.}\}$$

Thus, PDAs are in general nondeterministic because they may have a choice of transitions from any state. However, there are always only finitely many choices.

- An initial state  $q_0 \in Q$ ,
- An initial stack symbol  $Z_0 \in \Gamma$ ,
- A set of final states  $F \subseteq Q$ .

As an example we consider a PDA  $P$  that recognizes the language of even length palindromes over  $\Sigma = \{0, 1\}$ :  $L = \{ww^R \mid w \in \{0, 1\}^*\}$ . Intuitively, this PDA pushes the input symbols on the stack until it *guesses* that it is in the middle and then it compares the input with what is on the stack, popping of symbols from the stack as it goes. If it reaches the end of the input precisely at the time when the stack is empty, it accepts.

$$P_0 = (\{q_0, q_1, q_2\}, \{0, 1\}, \{0, 1, \#\}, \delta, q_0, \#, \{q_2\})$$

where  $\delta$  is given by the following equations:

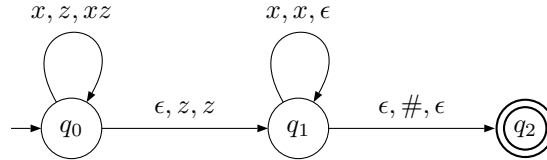
$$\begin{aligned}
\delta(q_0, 0, \#) &= \{(q_0, 0\#)\} \\
\delta(q_0, 1, \#) &= \{(q_0, 1\#)\} \\
\delta(q_0, 0, 0) &= \{(q_0, 00)\} \\
\delta(q_0, 1, 0) &= \{(q_0, 10)\} \\
\delta(q_0, 0, 1) &= \{(q_0, 01)\} \\
\delta(q_0, 1, 1) &= \{(q_0, 11)\} \\
\delta(q_0, \epsilon, \#) &= \{(q_1, \#)\} \\
\delta(q_0, \epsilon, 0) &= \{(q_1, 0)\} \\
\delta(q_0, \epsilon, 1) &= \{(q_1, 1)\} \\
\delta(q_1, 0, 0) &= \{(q_1, \epsilon)\} \\
\delta(q_1, 1, 1) &= \{(q_1, \epsilon)\} \\
\delta(q_1, \epsilon, \#) &= \{(q_2, \epsilon)\} \\
\delta(q, w, z) &= \emptyset \quad \text{everywhere else}
\end{aligned}$$

To save space we may abbreviate this by writing:

$$\begin{aligned}
\delta(q_0, x, z) &= \{(q_0, xz)\} \\
\delta(q_0, \epsilon, z) &= \{(q_1, z)\} \\
\delta(q_1, x, x) &= \{(q_1, \epsilon)\} \\
\delta(q_1, \epsilon, \#) &= \{(q_2, \epsilon)\} \\
\delta(q, x, z) &= \emptyset \quad \text{everywhere else}
\end{aligned}$$

where  $q \in Q, x \in \Sigma, z \in \Gamma$ . We obtain the previous table by expanding all the possibilities for  $q, x, z$ .

We draw the transition diagram of  $P$  by labelling each transition with a triple  $x, Z, \gamma$  with  $x \in \Sigma, Z \in \Gamma, \gamma \in \Gamma^*$ :



## 9.2 How does a PDA work?

At any time the state of the computation of a PDA  $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$  is given by:

- the state  $q \in Q$  the PDA is in,
- the input string  $w \in \Sigma^*$  that still has to be processed,
- the contents of the stack  $\gamma \in \Gamma^*$ .

Such a triple  $(q, w, \gamma) \in Q \times \Sigma^* \times \Gamma^*$  is called an Instantaneous Description (ID).

We define a relation  $\vdash_P \subseteq \text{ID} \times \text{ID}$  between IDs that describes how the PDA can change from one ID to the next one. Because PDAs in general are non-deterministic, this is a relation (not a function); i.e., there may be more than one possibility.

There are two possibilities for  $\vdash_P$ :

$$1. (q, xw, z\gamma) \vdash_P (q', w, \alpha\gamma) \text{ if } (q', \alpha) \in \delta(q, x, z)$$

$$2. (q, w, z\gamma) \vdash_P (q', w, \alpha\gamma) \text{ if } (q', \alpha) \in \delta(q, \epsilon, z)$$

In the first case the PDA reads an input symbol and consults the transition function  $\delta$  to calculate a possible new state  $q'$  and a sequence of stack symbols  $\alpha$  that replaces the current symbol on the top  $z$ .

In the second case the PDA ignores the input and silently moves into a new state and modifies the stack as above. The input is unchanged.

Consider the word 0110. What are possible sequences of IDs for  $P_0$  starting with  $(q_0, 0110, \#)$  ?

$$\begin{array}{ll} (q_0, 0110, \#) \vdash_{P_0} (q_0, 110, 0\#) & 1. \text{ with } (q_0, 0\#) \in \delta(q_0, 0, \#) \\ \vdash_{P_0} (q_0, 10, 10\#) & 1. \text{ with } (q_0, 10) \in \delta(q_0, 1, 0) \\ \vdash_{P_0} (q_1, 10, 10\#) & 2. \text{ with } (q_1, 1) \in \delta(q_0, \epsilon, 1) \\ \vdash_{P_0} (q_1, 0, 0\#) & 1. \text{ with } (q_1, \epsilon) \in \delta(q_1, 1, 1) \\ \vdash_{P_0} (q_1, \epsilon, \#) & 1. \text{ with } (q_1, \epsilon) \in \delta(q_1, 0, 0) \\ \vdash_{P_0} (q_2, \epsilon, \epsilon) & 2. \text{ with } (q_2, \epsilon) \in \delta(q_1, \epsilon, \#) \end{array}$$

We write  $(q, w, \gamma) \vdash_P^* (q', w', \gamma')$  if the PDA can move from  $(q, w, \gamma)$  to  $(q', w', \gamma')$  in a (possibly empty) sequence of moves. Above we have shown that

$$(q_0, 0110, \#) \vdash_{P_0}^* (q_2, \epsilon, \epsilon).$$

However, this is not the only possible sequence of IDs for this input. E.g. the PDA may just guess the middle wrong:

$$\begin{array}{ll} (q_0, 0110, \#) \vdash_{P_0} (q_0, 110, 0\#) & 1. \text{ with } (q_0, 0\#) \in \delta(q_0, 0, \#) \\ \vdash_{P_0} (q_1, 110, 0\#) & 2. \text{ with } (q_1, 0) \in \delta(q_0, \epsilon, 0) \end{array}$$

We have shown  $(q_0, 0110, \#) \vdash_{P_0}^* (q_1, 110, 0\#)$ . Here the PDA gets stuck as there is no state after  $(q_1, 110, 0\#)$ .

If we start with a word that is not in the language  $L$  (like 0011) then the automaton will always get stuck before reaching a final state.

### 9.3 The language of a PDA

There are two ways to define the language of a PDA  $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$  ( $L(P) \subseteq \Sigma^*$ ) because there are two notions of acceptance:

**Acceptance by final state**

$$L(P) = \{w \mid (q_0, w, Z_0) \vdash_P^* (q, \epsilon, \gamma) \wedge q \in F\}$$

That is the PDA accepts the word  $w$  if there is any sequence of IDs starting from  $(q_0, w, Z_0)$  and leading to  $(q, \epsilon, \gamma)$ , where  $q \in F$  is one of the final

states. Here it doesn't play a role what the contents of the stack are at the end.

In our example the PDA  $P_0$  would accept 0110 because  $(q_0, 0110, \#) \xrightarrow{*}_{P_0} (q_2, \epsilon, \epsilon)$  and  $q_2 \in F$ . Hence we conclude  $0110 \in L(P_0)$ .

On the other hand, because there is no successful sequence of IDs starting with  $(q_0, 0011, \#)$ , we know that  $0011 \notin L(P_0)$ .

### Acceptance by empty stack

$$L(P) = \{w \mid (q_0, w, Z_0) \xrightarrow{*}_P (q, \epsilon, \epsilon)\}$$

That is the PDA accepts the word  $w$  if there is any sequence of IDs starting from  $(q_0, w, Z_0)$  and leading to  $(q, \epsilon, \epsilon)$ , in this case the final state plays no role.

If we specify a PDA for acceptance by empty stack we will leave out the set of final states  $F$  and just use  $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0)$ .

Our example automaton  $P_0$  also works if we leave out  $F$  and use acceptance by empty stack.

We can always turn a PDA that uses one acceptance method into one that uses the other. Hence, both acceptance criteria specify the same class of languages.

## 9.4 Deterministic PDAs

We have introduced PDAs as nondeterministic machines that may have several alternatives how to continue. We now define Deterministic Pushdown Automata (DPDA) as those that never have a choice.

To be precise we say that a PDA  $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$  is deterministic (is a DPDA) iff

$$|\delta(q, x, z)| + |\delta(q, \epsilon, z)| \leq 1$$

Remember, that  $|X|$  stands for the number of elements in a finite set  $X$ .

That is: a DPDA may get stuck but it has never any choice.

In our example the automaton  $P_0$  is not deterministic, e.g. we have  $\delta(q_0, 0, \#) = \{(q_0, 0\#)\}$  and  $\delta(q_0, \epsilon, \#) = \{(q_1, \#)\}$  and hence  $|\delta(q_0, 0, \#)| + |\delta(q_0, \epsilon, \#)| = 2$ .

Unlike the situation for finite automata, there is in general no way to translate a nondeterministic PDA into a deterministic one. Indeed, there is no DPDA that recognizes the language  $L$ ! **Nondeterministic PDAs are more powerful than deterministic PDAs.**

However, we can define a similar language  $L'$  over  $\Sigma = \{0, 1, \$\}$  that can be recognized by a deterministic PDA:

$$L' = \{w\$w^R \mid w \in \{0, 1\}^*\}$$

That is  $L'$  contains palindroms with a marker  $\$$  in the middle, e.g.  $01\$10 \in L'$ .

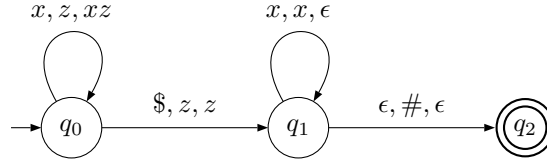
We define a DPDA  $P'$  for  $L'$ :

$$P' = (\{q_0, q_1, q_2\}, \{0, 1, \$\}, \{0, 1, \#\}, \delta', q_0, \#, \{q_2\})$$

where  $\delta'$  is given by:

$$\begin{aligned}
\delta'(q_0, x, z) &= \{(q_0, xz)\} & x \in \{0, 1\}, z \in \{0, 1, \#\} \\
\delta'(q_0, \$, z) &= \{(q_1, z)\} & z \in \{0, 1, \#\} \\
\delta'(q_1, x, x) &= \{(q_1, \epsilon)\} & x \in \{0, 1\} \\
\delta'(q_1, \epsilon, \#) &= \{(q_2, \epsilon)\} \\
\delta'(q, x, z) &= \emptyset & \text{everywhere else}
\end{aligned}$$

The transition graph is:



We can check that this automaton is deterministic. In particular the 3rd and 4th line cannot overlap because  $\#$  is not an input symbol.

In contrast to PDAs in general, the two acceptance methods are not equivalent for DPDAs: acceptance by final state makes it possible to define a bigger class of languages. We will consequently always use acceptance by final state for DPDAs in the following.

## 9.5 Context-free grammars and push-down automata

**Theorem 9.1** *For a language  $L \subseteq \Sigma^*$  the following two statements are equivalent:*

1.  $L$  is given by a CFG  $G$ ,  $L = L(G)$ .
2.  $L$  is the language of a PDA  $P$ ,  $L = L(P)$ .

To summarize: Context-Free Languages (CFLs) can be described by a Context-Free Grammar (CFG) and can be processed by a pushdown automaton.

We will only show how to construct a PDA from a grammar - the other direction is shown in [HMU01] (6.3.2, pp. 241).

Given a CFG  $G = (V, \Sigma, P, S)$ , we define a PDA

$$P(G) = (\{q_0\}, \Sigma, V \cup \Sigma, \delta, q_0, S)$$

where  $\delta$  is defined as follows:

$$\begin{aligned}
\delta(q_0, \epsilon, A) &= \{(q_0, \alpha) \mid A \rightarrow \alpha \in P\} & \text{for all } A \in V \\
\delta(q_0, a, a) &= \{(q_0, \epsilon)\} & \text{for all } a \in \Sigma.
\end{aligned}$$

We haven't given a set of final states because we use acceptance by empty stack. Yes, we use only one state!

Take as an example

$$G = (\{E, T, F\}, \{(\,, \,), \mathbf{a}, +, *\}, E, P)$$

where the set of productions  $P$  are given by

$$\begin{aligned}
E &\rightarrow T \mid E + T \\
T &\rightarrow F \mid T * F \\
F &\rightarrow \mathbf{a} \mid ( E )
\end{aligned}$$

we define

$$P(G) = (\{q_0\}, \{ (, ), \mathbf{a}, +, * \}, \{E, T, F, (, ), \mathbf{a}, +, * \}, \delta, q_0, E)$$

where

$$\begin{aligned} \delta(q_0, \epsilon, E) &= \{(q_0, T), (q_0, E+T)\} \\ \delta(q_0, \epsilon, T) &= \{(q_0, F), (q_0, T*F)\} \\ \delta(q_0, \epsilon, F) &= \{(q_0, \mathbf{a}), (q_0, (E))\} \\ \delta(q_0, (, ( &= \{(q_0, \epsilon)\} \\ \delta(q_0, ), ) &= \{(q_0, \epsilon)\} \\ \delta(q_0, \mathbf{a}, \mathbf{a}) &= \{(q_0, \epsilon)\} \\ \delta(q_0, +, +) &= \{(q_0, \epsilon)\} \\ \delta(q_0, *, *) &= \{(q_0, \epsilon)\} \\ \delta(q, x, z) &= \emptyset \quad \text{everywhere else} \end{aligned}$$

How does the  $P(G)$  accept  $\mathbf{a}+(\mathbf{a}*\mathbf{a})$ ?

$$\begin{aligned} (q_0, \mathbf{a}+(\mathbf{a}*\mathbf{a}), E) &\vdash (q_0, \mathbf{a}+(\mathbf{a}*\mathbf{a}), E+T) \\ &\vdash (q_0, \mathbf{a}+(\mathbf{a}*\mathbf{a}), T+T) \\ &\vdash (q_0, \mathbf{a}+(\mathbf{a}*\mathbf{a}), F+T) \\ &\vdash (q_0, \mathbf{a}+(\mathbf{a}*\mathbf{a}), \mathbf{a}+T) \\ &\vdash (q_0, +(\mathbf{a}*\mathbf{a}), +T) \\ &\vdash (q_0, (\mathbf{a}*\mathbf{a}), T) \\ &\vdash (q_0, (\mathbf{a}*\mathbf{a}), F) \\ &\vdash (q_0, (\mathbf{a}*\mathbf{a}), (E)) \\ &\vdash (q_0, \mathbf{a}*\mathbf{a}, E) \\ &\vdash (q_0, \mathbf{a}*\mathbf{a}, T) \\ &\vdash (q_0, \mathbf{a}*\mathbf{a}, T*F) \\ &\vdash (q_0, \mathbf{a}*\mathbf{a}, F*F) \\ &\vdash (q_0, \mathbf{a}*\mathbf{a}, \mathbf{a}*F) \\ &\vdash (q_0, *\mathbf{a}, *F) \\ &\vdash (q_0, \mathbf{a}, F) \\ &\vdash (q_0, \mathbf{a}, \mathbf{a}) \\ &\vdash (q_0, ), ) \\ &\vdash (q_0, \epsilon, \epsilon) \end{aligned}$$

Hence  $\mathbf{a}+(\mathbf{a}*\mathbf{a}) \in L(P(G))$ .

This above example illustrates the general idea:

$$\begin{aligned} w \in L(G) &\iff S \xRightarrow{*} w \\ &\iff (q_0, w, S) \vdash^* (q_0, \epsilon, \epsilon) \\ &\iff w \in L(P(G)) \end{aligned}$$

The automaton we have constructed is very nondeterministic: Whenever we have a choice between different rules the automaton may silently choose one of the alternatives.



## 10 Recursive-Descent Parsing

### 10.1 What is parsing?

According to Merriam-Webster OnLine<sup>20</sup> *parse* means:

To resolve (as a sentence) into component parts of speech and describe them grammatically.

In a computer science context, we take this to mean answering whether or not

$$w \in L(G)$$

for a CFG  $G$  by analysing the structure of  $w$  according to  $G$ ; i.e. to *recognise* the language generated by a grammar  $G$ .

A *parser* is a program that carries out parsing. For a CFG, this amounts to a realisation of a PDA. For most practical applications, a parser will also return a structured representation of a word  $w \in L(G)$ . This could be the derivation (or *parse*) *tree* for the word, or, more commonly, a simplified version of this, a so called *Abstract Syntax Tree*. Further, for a word not in the language,  $w \notin L(G)$ , a practical parser would normally provide an error message explaining why. An important application of parsers is in the front-end of compilers and interpreters where they turn a text-based representation of a program into a structured representation for further analysis and translation or execution.

In this section we study how to systematically construct a parser from a given CFG using the recursive-decent parsing method. To make the discussion concrete, we implement the parsers we develop in Haskell. Thus you can easily try out the examples in this section yourself using a Haskell system such as GHCi, and it is recommended that you do so! That said, we are essentially just using a small fragment of Haskell as a notation for writing simple functions, so it should not be difficult to follow the development even if you are not very familiar with Haskell. There are plenty of Haskell resources on-line, both for learning and for downloading and installing Haskell systems<sup>21</sup>.

### 10.2 Parsing strategies

There are two basic strategies for parsing: *top-down* and *bottom up*.

- A top-down parser attempts to carry out a derivation matching the input starting from the start symbol; i.e., it constructs the parse tree for the input *from the root downwards* in preorder.
- A bottom-up parser tries to construct the parse tree *from the leaves upwards* by using the productions “backwards”.

Top-down parsing is, in essence, what we have been doing so far whenever we have derived some specific word in a language from the start symbol of a grammar generating that language. For example, consider the grammar:

$$S \rightarrow aSa \mid bSb \mid a \mid b$$

---

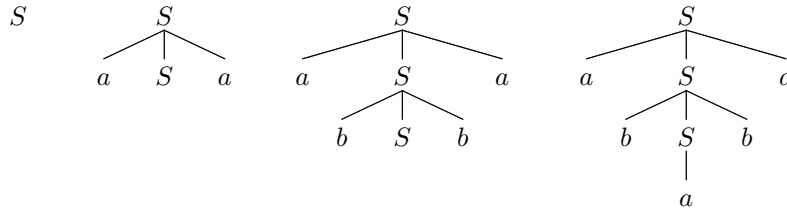
<sup>20</sup><http://www.webster.com>

<sup>21</sup><http://www.haskell.org>, <http://learnyouahaskell.com>

If given a string *ababa*, a top-down parser for this grammar would try to derive this string from the start symbol *S* and thus proceed as:

$$S \Rightarrow aSa \Rightarrow abSba \Rightarrow ababa$$

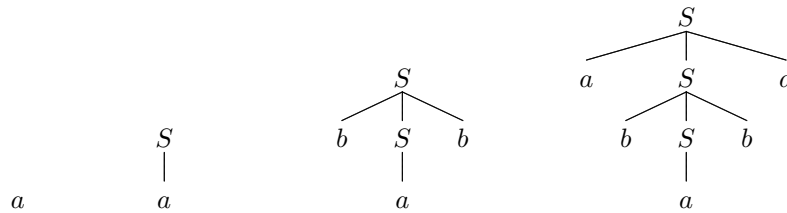
If we draw the derivation tree after each derivation step, we see how the tree gets constructed from the root downwards:



In contrast, a bottom-up parser would start from the leaves, and step by step group them together by applying productions in reverse:

$$ababa \Leftarrow abSba \Leftarrow aSa \Leftarrow S$$

This is a *rightmost derivation in reverse*. The derivation tree thus gets constructed from the bottom upwards:



The key difficulty of bottom-up parsing is to decide when to reduce. For example, in this case, how does the parser know not to reduce neither the first input symbol *a* nor the second input symbol *b* to *S*, but wait until it sees the middle *a* and only then do the first reduction step? Such questions are answered by *LR parsing theory*<sup>22</sup> We will not consider LR parsing further here, except noting that it is covered in more depth in the Compilers module [Nil16] that builds on much of the material in this module. Instead we turn our attention to recursive-decent parsing, which is a type of top-down parsing.

### 10.3 Basics of recursive-descent parsing

*Recursive-descent parsing* is a way to implement top-down parsing. We are just going to focus on the language recognition problem:  $w \in L(G)$ ? This suggests the following type for the parser:

```
parser :: [Token] -> Bool
```

*Token* is “compiler speak” for (input) symbol; i.e., an element of the alphabet.

Consider a typical production in some CFG *G*:

$$S \rightarrow AB$$

<sup>22</sup>[https://en.wikipedia.org/wiki/LR\\_parser](https://en.wikipedia.org/wiki/LR_parser)

Let  $L(X)$  be the language  $\{w \in T^* \mid X \xrightarrow[G]{*} w\}$ ,  $X \in N$ . Note that

$$w \in L(S) \Leftarrow \exists w_1, w_2 . \begin{array}{l} w = w_1 w_2 \\ \wedge w_1 \in L(A) \\ \wedge w_2 \in L(B) \end{array}$$

That is, given a parser for  $L(A)$  and a parser for  $L(B)$ , we can construct a parser for  $L(S)$  by asking the first parser if a prefix  $w_1$  of  $w$  belongs to  $L(A)$ , and then asking the other parser if the remaining suffix  $w_2$  of  $w$  belongs to  $L(B)$ . If the answer to both questions is yes, then  $w$  belongs to  $L(S)$ .

However, we need to find the right way to divide the input word  $w$ ! In general, there are  $|w| + 1$  possibilities. We could, of course, blindly try them all. But as the prefix and suffix recursively also have to be split in all possible ways, and so on until we get down to individual alphabet symbols, it is clear that this approach would lead to a combinatorial explosion that would render such a parser useless for all but very short words.

Instead we need to let the input guide the search. To that end, we initially adopt the following idea:

- Each parser tries to derive a *prefix* of the input according to the productions for the nonterminal
- Each parser returns the remaining *suffix* if successful, allowing this to be passed to the next parser for analysis.

This gives us the following refined type for parsers:

```
parseX :: [Token] -> Maybe [Token]
```

Recall that `Maybe` is Haskell's option type:

```
data Maybe a = Nothing | Just a
```

Of course, we should be a little suspicious: There could be *more* than one prefix derivable from a non-terminal, and if so, how can we then know which *one* to pick? Picking the *wrong* prefix might make it impossible to derive the suffix from the non-terminal that follows. We will return to these points later.

Now we can construct a parser for  $L(S)$

$$S \rightarrow AB$$

in terms of parsers for  $L(A)$  and  $L(B)$ :

```
parseS :: [Token] -> Maybe [Token]
parseS ts =
  case parseA ts of
    Nothing -> Nothing
    Just ts' ->
      case parseB ts' of
        Nothing -> Nothing
        Just ts'' -> Just ts''
```

Note that the case analysis on the result of `parseB ts'` simply passes on the result unchanged. Thus we can simplify the code to:

```

parseS :: [Token] -> Maybe [Token]
parseS ts =
  case parseA ts of
    Nothing -> Nothing
    Just ts' -> parseB ts'

```

This approach is called recursive-descent parsing because the parse functions (usually) end up being (mutually) recursive. What does this have to do with realising a PDA? Fundamental to the implementation of a recursive computation is a that keeps track of the *state* of the computation and allows for *subcomputations* (to any depth). In a language that supports recursive functions and procedures, the stack is usually not explicitly visible, but internally, it is a central datastructure. Thus, a recursive-descent parser is a kind of PDA.

Let us develop this example a little further into code that can be executed. First, for simplicity, let us pick the type `Char` for token:

```

type Token = Char

```

Recall that (basic) strings are just lists of characters in Haskell; that is, `String = [Char] = [Token]`. Thus, a string literal like `"abcd"` is just a shorthand notation for the list of characters `['a','b','c','d']`.

Now, suppose the productions for  $A$  and  $B$  are the following

$$\begin{aligned}
 A &\rightarrow a \\
 B &\rightarrow b
 \end{aligned}$$

Thus, in this case, it is clear that if the parsing function `parseA` for the non-terminal  $A$  sees input starting with an  $a$ , then that  $a$  is the desired prefix, and whatever remains of the input is the suffix that should be returned as part if the indication of having been able to successfully derive a prefix of the input from the nonterminal in question. Otherwise, if the input does not start with an  $a$ , it is equally clear that it is not possible to derive any prefix of the input from the nonterminal  $A$ , and the parsing function must thus indicate failure. In essence, this is how the input is used to guide the search for the right prefix.

We can implement `parseA` in Haskell, using pattern matching, as follows:

```

parseA :: [Token] -> Maybe [Token]
parseA ('a' : ts) = Just ts
parseA _         = Nothing

```

The case and code for the parsing function `parseB` for the nonterminal  $B$  is of course analogous:

```

parseB :: [Token] -> Maybe [Token]
parseB ('b' : ts) = Just ts
parseB _         = Nothing

```

Now we can evaluate `parseA`, `parseB`, and `parseS` on `"abcd"` with the following results:

```

parseA "abcd" => Just "bcd"
parseB "abcd" => Nothing
parseS "abcd" => Just "cd"

```

This tells us that a prefix of  $abcd$  can be derived from  $A$ , leaving a remaining suffix  $bcd$ , that no prefix of  $abcd$  can be derived from  $B$ , but that a prefix of  $abcd$  also can be derived from  $S$ , leaving a suffix  $cd$ , as we would expect.

## 10.4 Handling choice

Of course, there are usually more than one one production for a nonterminal. Thus need a way to handle *choice*, as in

$$S \rightarrow AB \mid CD$$

We are first going to consider the case when the choice is obvious, as in

$$S \rightarrow aAB \mid cCD$$

That is, we assume it is manifest from the grammar that we can choose between productions with a one-symbol *lookahead*.

As an example, let us construct a parser for the grammar:

$$\begin{aligned} S &\rightarrow aA \mid bBA \\ A &\rightarrow aA \mid \epsilon \\ B &\rightarrow bB \mid \epsilon \end{aligned}$$

We are going to need one parsing function for each non-terminal:

- `parseS :: [Token] -> Maybe [Token]`
- `parseA :: [Token] -> Maybe [Token]`
- `parseB :: [Token] -> Maybe [Token]`

We again take `type Token = Char` for simplicity.

Code for `parseS`. Note how the pattern matching makes use of a one-symbol lookahead to chose between the two productions for  $S$ :

```
parseS :: [Token] -> Maybe [Token]
parseS ('a' : ts) =
    parseA ts
parseS ('b' : ts) =
    case parseB ts of
        Nothing -> Nothing
        Just ts' -> parseA ts'
parseS _ = Nothing
```

The code for `parseA` and `parseB` similarly make use of the one-symbol lookahead to chose between productions, but this time, because  $A \Rightarrow \epsilon$  and  $B \Rightarrow \epsilon$ , it is *not* a syntax error if the next token is not, respectively,  $a$  and  $b$ . Thus both functions can succeed *without* consuming any input:

```
parseA :: [Token] -> Maybe [Token]
parseA ('a' : ts) = parseA ts
parseA ts         = Just ts

parseB :: [Token] -> Maybe [Token]
parseB ('b' : ts) = parseB ts
parseB ts         = Just ts
```

Now consider a more challenging scenario:

$$\begin{aligned} S &\rightarrow aA \mid aBA \\ A &\rightarrow aA \mid \epsilon \\ B &\rightarrow bB \mid \epsilon \end{aligned}$$

In the parsing function `parseS` for nonterminal  $S$ , should `parseA` or `parseB` be called once `a` has been read?

We could try the alternatives in order; i.e., a limited form of *backtracking*:

```
parseS ('a' : ts) =
  case parseA ts of
    Just ts' -> Just ts'
    Nothing ->
      case parseB ts of
        Nothing -> Nothing
        Just ts' -> parseA ts'
```

Of course, the choice to try `parseA` first is arbitrary, a point we will revisit shortly.

Similarly, there are two alternatives for the nonterminal  $A$ . In fact, we already encountered this situation above, and as we did there, let us try to consume an input prefix (here  $a$ ) if possible, and only if that is not possible succeed without consuming any input:

```
parseA :: [Token] -> Maybe [Token]
parseA ('a' : ts) = parseA ts
parseA ts         = Just ts
```

The code for  $B$  is of course similar. This may seem like an obvious ordering: after all, if we opted to “try” without consuming any input first, then that would always succeed, and no other alternatives would ever be tried. Nevertheless, picking the order we did still amounts to an arbitrary choice, and in fact it is *not* always the right one.

The problem here is that limited backtracking is *not* an exhaustive search. For many grammars, there simply is no one order that *always* will work, meaning that a parser that nevertheless commits to one particular order is liable to get stuck in “blind alleys”.

Consider the following grammar

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aA \mid \epsilon \\ B &\rightarrow ab \end{aligned}$$

and corresponding parsing functions:

```
parseA ('a' : ts) = parseA ts
parseA ts         = Just ts

parseB ('a' : 'b' : ts) = Just ts
parseB ts               = Nothing
```

```

parseS ts =
  case parseA ts of
    Nothing  -> Nothing
    Just ts' -> parseB ts'

```

Will it work? Let us try it on *ab*. Clearly derivable from the grammar:

$$S \Rightarrow AB \Rightarrow B \Rightarrow ab$$

However, if we run our parser on "ab":

```

parseS "ab"  => Nothing

```

Our parser thus says “no”. Why? Because

```

parseA "ab"  => Just "b"

```

That is, the code for `parseA` committed to the choice  $A \rightarrow a$  too early and will never try  $A \rightarrow \epsilon$ . That was the wrong choice in this case and the parser got stuck in a “blind alley”.

Would it have been better to try  $A \rightarrow \epsilon$  first? Then the parser would work for the word *ab*, but it would still fail on other words that should be accepted, such as *aab*. To successfully parse that word, `parseA` must somehow consume the first *a* but not the second, and neither ordering of the productions for *A* will achieve that.

One principled approach addressing this dilemma is to try *all* alternatives; i.e., *full backtracking* (aka *list of successes*):

- Each parsing function returns a *list* of *all* possible suffixes. Type:

```

parseX :: [Token] -> [[Token]]

```

- Translate  $A \rightarrow \alpha \mid \beta$  into

```

parseA ts = parseAlpha ts ++ parseBeta ts

```

- An empty list indicates no possible parsing.

However:

- Full backtracking is computationally expensive.
- In error reporting, it becomes difficult to pinpoint the exact location of a syntax error: where exactly lies the problem if it only *after* an exhaustive search becomes apparent that there is no possible way to parse a word?

In section 10.6, we are going to look at another principled approach that avoids backtracking: *predictive parsing*. The price we have to pay is that the grammar must satisfy certain conditions. But at least we will know *statically, at construction time* if the parser is going to work or not. And if not, we can try to modify the grammar (without changing the language) until the prerequisite conditions are met. First, however, we will consider the problem of left-recursion and context-free grammars.

## 10.5 Recursive-descent parsing and left-recursion

Consider the grammar

$$A \rightarrow Aa \mid \epsilon$$

and the corresponding recursive-descent parsing function:

```
parseA :: [Token] -> Maybe [Token]
parseA ts =
  case parseA ts of
    Just ('a' : ts') -> Just ts'
    -                 -> Just ts
```

Any problem? Yes, because the function calls itself without consuming any input, it will loop forever.

The problem here is that the grammar is left-recursive. Recall that a this means that there is a derivation  $A \xRightarrow{+} A\alpha$  for some nonterminal  $A$ . As each derivation step corresponds to one parsing function calling another, it is clear that recursive-descent parsing functions derived from a left-recursive grammar will end up looping forever as soon as one of the parsing functions for a left-recursive non-terminal is invoked because no input is consumed before the same function is entered again, directly or indirectly.

Recursive-descent parsers thus *cannot*<sup>23</sup> deal with *left-recursive* grammars. The standard way of resolving this is to transform a left-recursive grammar into an equivalent grammar that is not left recursive as described in section 8.6, and then deriving the parser from the non-left-recursive version of the grammar.

## 10.6 Predictive parsing

In a recursive-decent parsing setting, we want a parsing function to be successful *exactly* when a prefix of the input *can* be derived from the corresponding nonterminal. This can be achieved by:

- Adopting a suitable parsing strategy, specifically regarding how to handle *choice* between two or more productions for one nonterminal.
- Impose *restrictions* on the grammar to ensure success of the chosen parsing strategy.

*Predictive parsing* is when all parsing decisions can be made based on a *lookahead* of limited length, typically one symbol. We have already seen cases where predictive parsing clearly is possible; for example, this is manifestly the case when the right-hand side of each possible production starts with a distinct terminal, as here:

$$S \rightarrow aB \mid cD$$

But we also saw that the choice is not always this obvious, and that if we then make arbitrary choices regarding which order in which to try productions, the resulting parser is likely to be flawed. In the following, we are going to

---

<sup>23</sup>At least not when implemented in the standard way described here. It is possible, by keeping track of more contextual information, to detect when no progress is being made and limit the recursion depth at that point. See [https://en.wikipedia.org/wiki/Top-down\\_parsing](https://en.wikipedia.org/wiki/Top-down_parsing).



look into *exactly* when the next input symbol suffices to make all choices. As a consequence, if we are faced with a grammar where a one-symbol lookahead is not enough, we will know this, and we can take corrective action, such as trying to transform the grammar in a way that will resolve the problem.

Before we start, let us just give an example that illustrates that a one-symbol lookahead can be enough even if the RHSs start with nonterminals:

$$\begin{aligned} S &\rightarrow AB \mid CD \\ A &\rightarrow a \mid b \\ C &\rightarrow c \mid d \end{aligned}$$

Here, if the input starts with an  $a$  or  $b$ , we should clearly attempt to parse by the production  $S \rightarrow AB$ , and if it starts with a  $c$  or a  $d$ , we should attempt to parse by  $S \rightarrow CD$ . This suggests that the key is going to be an *analysis* of the grammar: for each nonterminal, we need to know what symbols that may start words derived from that nonterminal.

More generally, consider productions for a nonterminal  $X$

$$X \rightarrow \alpha \mid \beta$$

and the corresponding parsing code:

```
parseX (t : ts) =
  | t ??      -> parse α
  | t ??      -> parse β
  | otherwise -> Nothing
```

The question is, what should the conditions be on the lookahead symbol  $t$ , here indicated by  $??$ , to decide whether or not to parse by the production corresponding to each case?

The idea of predictive parsing is this:

- Compute the *set* of terminal symbols that can *start* strings derived from each alternative, the *first set*.
- If there is a choice between two or more alternatives, insist that the first sets for those are *disjoint* (a grammar restriction).
- The right choice can now be made simply by determining to which alternative's first set the next input symbol belongs.

We can now refine the code to:

```
parseX (t : ts) =
  | t ∈ first(α) -> parse α
  | t ∈ first(β) -> parse β
  | otherwise   -> Nothing
```

But the situation could be a bit more involved as it sometimes is possible to derive the empty word from a nonterminal, and the empty word does of course not begin with any symbol at all. For a concrete example, consider again  $X \rightarrow \alpha \mid \beta$ , and suppose it can be the case that  $\beta \xRightarrow{*} \epsilon$ .

Clearly, the next input symbol could in this case be a terminal that can follow a string derivable from  $X$ , meaning we need to refine the parsing code further:

```

parseX (t : ts) =
  | t ∈ first(α)           -> parse α
  | t ∈ first(β) ∪ follow(X) -> parse β
  | otherwise              -> Nothing

```

Of course, the branches must be mutually exclusive! Otherwise a one-symbol lookahead is not enough to decide which choice to make.

### 10.6.1 First and follow sets

We will now develop these ideas in more detail. The presentation roughly follows “The Dragon Book” [ASU86]. For a CFG  $G = (N, T, P, S)$ :

$$\begin{aligned}
 \text{first}(\alpha) &= \{a \in T \mid \alpha \xrightarrow{*}_G a\beta\} \\
 \text{follow}(A) &= \{a \in T \mid S \xrightarrow{*}_G \alpha A a \beta\} \\
 &\quad \cup \{\$ \mid S \xrightarrow{*}_G \alpha A\}
 \end{aligned}$$

where  $\alpha, \beta \in (N \cup T)^*$ ,  $A \in N$ , and where  $\$$  is a special “end of input” marker.

To illustrate these definitions, consider the grammar:

$$\begin{array}{ll}
 S \rightarrow ABC & B \rightarrow b \mid \epsilon \\
 A \rightarrow aA \mid \epsilon & C \rightarrow c \mid d
 \end{array}$$

First sets:

$$\begin{aligned}
 \text{first}(C) &= \{c, d\} \\
 \text{first}(B) &= \{b\} \\
 \text{first}(A) &= \{a\} \\
 \text{first}(S) &= \text{first}(ABC) \\
 &= [\text{because } A \xrightarrow{*} \epsilon \text{ and } B \xrightarrow{*} \epsilon] \\
 &\quad \text{first}(A) \cup \text{first}(B) \cup \text{first}(C) \\
 &= \{a, b, c, d\}
 \end{aligned}$$

Follow sets:

$$\begin{aligned}
 \text{follow}(C) &= \{\$ \} \\
 \text{follow}(B) &= \text{first}(C) = \{c, d\} \\
 \text{follow}(A) &= [\text{because } B \xrightarrow{*} \epsilon] \\
 &\quad \text{first}(B) \cup \text{first}(C) \\
 &= \{b, c, d\}
 \end{aligned}$$

### 10.6.2 LL(1) grammars

Now consider *all* productions for a nonterminal  $A$  in some grammar:

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

In the parsing function for  $A$ , on input symbol  $t$ , we should parse according to  $\alpha_i$  if

- $t \in \text{first}(\alpha_i)$ .
- $t \in \text{follow}(A)$ , if  $\alpha_i \xrightarrow{*} \epsilon$

Thus, if:

- $\text{first}(\alpha_i) \cap \text{first}(\alpha_j) = \emptyset$  for  $1 \leq i < j \leq n$ , and
- if  $\alpha_i \xrightarrow{*} \epsilon$  for some  $i$ , then, for all  $1 \leq j \leq n$ ,  $j \neq i$ ,
  - $\alpha_j \not\xrightarrow{*} \epsilon$ , and
  - $\text{follow}(A) \cap \text{first}(\alpha_j) = \emptyset$

then it is always clear what to do! A grammar satisfying these conditions is said to be an *LL(1)* grammar.

### 10.6.3 Nullable nonterminals

In order to compute the first and follow sets for a grammar  $G = (N, T, P, S)$ , we first need to know all nonterminals  $A \in N$  such that  $A \xrightarrow{*} \epsilon$ ; i.e. the set  $N_\epsilon \subseteq N$  of *nullable* nonterminals.

Let  $\text{syms}(\alpha)$  denote the *set* of symbols in a string  $\alpha$ :

$$\begin{aligned} \text{syms} &\in (N \cup T)^* \rightarrow \mathcal{P}(N \cup T) \\ \text{syms}(\epsilon) &= \emptyset \\ \text{syms}(X\alpha) &= \{X\} \cup \text{syms}(\alpha) \end{aligned}$$

The set  $N_\epsilon$  is the *smallest* solution to the equation

$$N_\epsilon = \{A \mid A \rightarrow \alpha \in P \wedge \forall X \in \text{syms}(\alpha) . X \in N_\epsilon\}$$

Note that  $A \in N_\epsilon$  if  $A \rightarrow \epsilon \in P$  because  $\text{syms}(\epsilon) = \emptyset$  and  $\forall X \in \emptyset . \dots$  is trivially true. Also note that we really need to look for the *smallest* solution. For example, consider a grammar  $S \rightarrow SS \mid a$ .  $N_\epsilon = \{S\}$  is clearly a solution to the equation defining  $N_\epsilon$  for this grammar, but  $S$  is also clearly not nullable. That is because  $N_\epsilon = \{S\}$  is *not* the smallest solution. The smallest solution in this case is  $N_\epsilon = \emptyset$ ; i.e. there are no nullable nonterminals.

We can now define a predicate nullable on *strings* of grammar symbols:

$$\begin{aligned} \text{nullable} &\in (N \cup T)^* \rightarrow \text{Bool} \\ \text{nullable}(\epsilon) &= \text{true} \\ \text{nullable}(X\alpha) &= X \in N_\epsilon \wedge \text{nullable}(\alpha) \end{aligned}$$

The equation for  $N_\epsilon$  can be solved iteratively as follows:

1. Initialize  $N_\epsilon$  to  $\{A \mid A \rightarrow \epsilon \in P\}$ .
2. If there is a production  $A \rightarrow \alpha$  such that  $\forall X \in \text{syms}(\alpha) . X \in N_\epsilon$ , then add  $A$  to  $N_\epsilon$ .
3. Repeat step 2 until no further nullable nonterminals can be found.

Consider the following grammar:

$$\begin{array}{ll} S \rightarrow ABC \mid AB & B \rightarrow b \mid \epsilon \\ A \rightarrow aA \mid BB & C \rightarrow c \mid d \end{array}$$

- Because  $B \rightarrow \epsilon$  is a production,  $B \in N_\epsilon$ .
- Because  $A \rightarrow BB$  is a production and  $B \in N_\epsilon$ , additionally  $A \in N_\epsilon$ .
- Because  $S \rightarrow AB$  is a production, and  $A, B \in N_\epsilon$ , additionally  $S \in N_\epsilon$ .
- No more production with nullable RHSs. The set of nullable symbols  $N_\epsilon = \{S, A, B\}$ .

#### 10.6.4 Computing first sets

For a CFG  $G = (N, T, P, S)$ , the sets  $\text{first}(A)$  for  $A \in N$  are the smallest sets satisfying:

$$\begin{aligned} \text{first}(A) &\subseteq T \\ \text{first}(A) &= \bigcup_{A \rightarrow \alpha \in P} \text{first}(\alpha) \end{aligned}$$

For strings,  $\text{first}$  is defined as (note the *overloaded* notation):

$$\begin{aligned} \text{first} &\in (N \cup T)^* \rightarrow \mathcal{P}(T) \\ \text{first}(\epsilon) &= \emptyset \\ \text{first}(a\alpha) &= \{a\} \\ \text{first}(A\alpha) &= \text{first}(A) \cup \begin{cases} \text{first}(\alpha), & \text{if } A \in N_\epsilon \\ \emptyset, & \text{if } A \notin N_\epsilon \end{cases} \end{aligned}$$

where  $a \in T$ ,  $A \in N$ , and  $\alpha \in (N \cup T)^*$ .

The solutions can often be obtained directly by expanding out all definitions. If necessary, the equations can be solved by iteration in a similar way to how  $N_\epsilon$  is computed. Note that the smallest solution to a set equation of the type  $X = X \cup Y$  when there are no other constraints on  $X$  is simply  $X = Y$ .

Consider (again):

$$\begin{array}{ll} S \rightarrow ABC & B \rightarrow b \mid \epsilon \\ A \rightarrow aA \mid \epsilon & C \rightarrow c \mid d \end{array}$$

First compute the nullable nonterminals:  $N_\epsilon = \{A, B\}$ . Then the first sets:

$$\begin{aligned}
\text{first}(A) &= \text{first}(aA) \cup \text{first}(\epsilon) \\
&= \{a\} \cup \emptyset = \{a\} \\
\text{first}(B) &= \text{first}(b) \cup \text{first}(\epsilon) \\
&= \{b\} \cup \emptyset = \{b\} \\
\text{first}(C) &= \text{first}(c) \cup \text{first}(d) \\
&= \{c\} \cup \{d\} = \{c, d\} \\
\text{first}(S) &= \text{first}(ABC) \\
&= [A \in N_\epsilon] \\
&\quad \text{first}(A) \cup \text{first}(BC) \\
&= [B \in N_\epsilon \wedge C \notin N_\epsilon] \\
&\quad \text{first}(A) \cup \text{first}(B) \cup \text{first}(C) \cup \emptyset \\
&= \{a\} \cup \{b\} \cup \{c, d\} = \{a, b, c, d\}
\end{aligned}$$

### 10.6.5 Computing follow sets

For a CFG  $G = (N, T, P, S)$ , the sets  $\text{follow}(A)$  are the smallest sets satisfying:

- $\{\$ \} \subseteq \text{follow}(S)$
- If  $A \rightarrow \alpha B \beta \in P$ , then  $\text{first}(\beta) \subseteq \text{follow}(B)$
- If  $A \rightarrow \alpha B \beta \in P$ , and  $\text{nullable}(\beta)$  then  $\text{follow}(A) \subseteq \text{follow}(B)$

$A, B \in N$ , and  $\alpha, \beta \in (N \cup T)^*$ .

(It is assumed that there are no *useless* symbols; i.e., all symbols can appear in the derivation of some sentence.)

Here is our example grammar again:

$$\begin{array}{ll}
S \rightarrow ABC & B \rightarrow b \mid \epsilon \\
A \rightarrow aA \mid \epsilon & C \rightarrow c \mid d
\end{array}$$

Constraints for  $\text{follow}(S)$ :

$$\{\$ \} \subseteq \text{follow}(S)$$

Constraints for  $\text{follow}(A)$  (note:  $\neg \text{nullable}(BC)$ ):

$$\begin{aligned}
\text{first}(BC) &\subseteq \text{follow}(A) \\
\text{first}(\epsilon) &\subseteq \text{follow}(A) \\
\text{follow}(A) &\subseteq \text{follow}(A)
\end{aligned}$$

Constraints for  $\text{follow}(B)$  (note:  $\neg \text{nullable}(C)$ ):

$$\text{first}(C) \subseteq \text{follow}(B)$$

Constraints for  $\text{follow}(C)$  (note:  $\text{nullable}(\epsilon)$ ):

$$\begin{aligned}
\text{first}(\epsilon) &\subseteq \text{follow}(C) \\
\text{follow}(S) &\subseteq \text{follow}(C)
\end{aligned}$$

In general:

$$X \subseteq Z \wedge Y \subseteq Z \iff X \cup Y \subseteq Z$$

Also, constraints like  $\emptyset \subseteq X$  and  $X \subseteq X$  are trivially satisfied and can be omitted.

The constraints for our example can thus be written as:

$$\begin{aligned} \{\$ \} &\subseteq \text{follow}(S) \\ \text{first}(BC) \cup \text{first}(\epsilon) &\subseteq \text{follow}(A) \\ \text{first}(C) &\subseteq \text{follow}(B) \\ \text{first}(\epsilon) \cup \text{follow}(S) &\subseteq \text{follow}(C) \end{aligned}$$

Using

$$\begin{aligned} \text{first}(\epsilon) &= \emptyset \\ \text{first}(C) &= \{c, d\} \\ \text{first}(BC) &= \text{first}(B) \cup \text{first}(C) \cup \emptyset \\ &= \{b\} \cup \{c, d\} = \{b, c, d\} \end{aligned}$$

the constraints can be simplified further:

$$\begin{aligned} \{\$ \} &\subseteq \text{follow}(S) \\ \{b, c, d\} &\subseteq \text{follow}(A) \\ \{c, d\} &\subseteq \text{follow}(B) \\ \text{follow}(S) &\subseteq \text{follow}(C) \end{aligned}$$

Finally, looking for the smallest sets satisfying these constraints, we get:

$$\begin{aligned} \text{follow}(S) &= \{\$ \} \\ \text{follow}(A) &= \{b, c, d\} \\ \text{follow}(B) &= \{c, d\} \\ \text{follow}(C) &= \text{follow}(S) = \{\$ \} \end{aligned}$$

### 10.6.6 Implementing a predictive parser

Let us now implement a predictive parser for our sample grammar:

$$\begin{array}{ll} S \rightarrow ABC & B \rightarrow b \mid \epsilon \\ A \rightarrow aA \mid \epsilon & C \rightarrow c \mid d \end{array}$$

For this grammar, as per the calculations above:

- Nullable symbols:  $N_\epsilon = \{S, A, B\}$
- First sets:

$$\begin{aligned} \text{first}(S) &= \{a, b, c, d\} \\ \text{first}(A) &= \{a\} \\ \text{first}(B) &= \{b\} \\ \text{first}(C) &= \{c, d\} \end{aligned}$$

- Follow sets:

$$\begin{aligned}
\text{follow}(S) &= \{\$ \} \\
\text{follow}(A) &= \{b, c, d\} \\
\text{follow}(B) &= \{c, d\} \\
\text{follow}(C) &= \{\$ \}
\end{aligned}$$

Recall the “template code” for a parsing function for a pair of productions like  $X \rightarrow \alpha \mid \beta$ . If no RHS is nullable:

```

parseX (t : ts) =
  | t ∈ first(α) -> parse α
  | t ∈ first(β) -> parse β
  | otherwise   -> Nothing

```

If one RHS is nullable, say nullable( $\beta$ ):

```

parseX (t : ts) =
  | t ∈ first(α)                -> parse α
  | t ∈ first(β) ∪ follow(X) -> parse β
  | otherwise                   -> Nothing

```

We can now implement predictive parsing functions for the nonterminals  $S$ ,  $A$ ,  $B$ , and  $C$  as follows in pseudo Haskell. The function for  $S$  is straightforward as there is only one production for  $S$ , thus no choice:

```

parseS ts =
  case parseA ts of
    Just ts' ->
      case parseB ts' of
        Just ts'' ->
          parseC ts''
        Nothing ->
          Nothing
    Nothing ->
      Nothing

```

For  $A$ , there are two productions:  $A \rightarrow aA \mid \epsilon$ . Note that  $aA$  is trivially not nullable, while  $\epsilon$  trivially is nullable. We thus compute  $\text{first}(aA) = \{a\}$  and  $\text{first}(\epsilon) \cup \text{follow}(A) = \emptyset \cup \{b, c, d\} = \{b, c, d\}$ . The parsing function for  $A$  thus becomes:

```

parseA (t : ts) =
  | t ∈ {a}      -> parseA ts
  | t ∈ {b, c, d} -> Just (t : ts)
  | otherwise    -> Nothing

```

Note how the case for the  $\epsilon$ -production only does checking on the next input symbol, but does not consume it.

For  $B$ , there are also two productions:  $B \rightarrow b \mid \epsilon$ . Note that  $b$  is trivially not nullable, while  $\epsilon$  trivially is nullable. We thus compute  $\text{first}(b) = \{b\}$  and  $\text{first}(\epsilon) \cup \text{follow}(B) = \emptyset \cup \{c, d\} = \{c, d\}$ . The resulting parsing function for  $B$ :

```

parseB (t : ts) =
  | t ∈ {b}   -> Just ts
  | t ∈ {c,d} -> Just (t : ts)
  | otherwise -> Nothing

```

Note how the case for  $b$  checks that the next input indeed is a  $b$  and if so succeeds, consuming that one input symbol, while the  $\epsilon$ -production again only checks the next input symbol without consuming it.

Finally, the productions for  $C$  are  $C \rightarrow c \mid d$ , where both RHSs are trivially not nullable. We compute  $\text{first}(c) = \{c\}$  and  $\text{first}(d) = \{d\}$ . The parsing function for  $C$ :

```

parseC (t : ts) =
  | t ∈ {c}   -> Just ts
  | t ∈ {d}   -> Just ts
  | otherwise -> Nothing

```

This can of course be simplified a little, which is the case whenever there are multiple productions with the RHSs being a single terminal:

```

parseC (t : ts) =
  | t ∈ {c,d} -> Just ts
  | otherwise -> Nothing

```

### 10.6.7 LL(1), left-recursion, and ambiguity

As we have seen, the LL(1) conditions impose a number of restrictions on a grammar. In particular, no left-recursive or ambiguous grammar can be LL(1)!

Let us prove that a left-recursive grammar cannot be LL(1). Recall that a grammar is left-recursive iff there exists  $A \in N$  and  $\alpha \in (N \cup T)^*$  such that  $A \xRightarrow{+} A\alpha$  (section 8.6). We can assume without loss of generality that there are no useless symbols and productions in the grammar as any useless productions can be removed from a grammar without changing the language (section 8.2). It thus follows that a derivation  $A \xRightarrow{+} w$ ,  $w \in T^*$ , must also exist.

Let us assume that all derivations are leftmost. Clearly,  $A\alpha \neq w$ , and thus there must have been a choice at some point differentiating these two derivations. That is, there must exist some  $B \in N$  for which there are at least two *distinct* productions  $B \rightarrow \beta_1 \mid \beta_2$  such that

$$A \xRightarrow{*} B\gamma \Rightarrow \beta_1\gamma \xRightarrow{*} A\alpha$$

and

$$A \xRightarrow{*} B\gamma \Rightarrow \beta_2\gamma \xRightarrow{*} w$$

Let us now observe that if there is a derivation  $\alpha \xRightarrow{*} \beta$ , then  $\text{first}(\alpha) \supseteq \text{first}(\beta)$ . Let us also observe that if  $\neg \text{nullable}(\alpha)$ , then  $\text{first}(\alpha\beta) = \text{first}(\alpha)$ , and if  $\text{nullable}(\alpha)$ , then  $\text{first}(\alpha\beta) = \text{first}(\alpha) \cup \text{first}(\beta)$ . (These should really be proved as auxiliary lemmas, but they are fairly obvious.)

Now let us consider  $\beta_1$  and  $\beta_2$ . If both  $\text{nullable}(\beta_1)$  and  $\text{nullable}(\beta_2)$ , then that is an immediate violation of the LL(1) conditions, so we need not consider that case further. Moreover, we can assume  $w \neq \epsilon$ : if the only terminal string derivable from  $A$  is  $\epsilon$ , then, under the assumption of no useless productions, it



must be the case that both  $\text{nullable}(\beta_1)$  and  $\text{nullable}(\beta_2)$  which again violates the LL(1) conditions. Thus, because  $A \xRightarrow{+} w$  and  $w \neq \epsilon$ , we have:

$$\text{first}(A) \supseteq \text{first}(w) \neq \emptyset$$

Suppose  $\neg\text{nullable}(\beta_1)$  and  $\neg\text{nullable}(\beta_2)$ . Because  $\neg\text{nullable}(\beta_1)$ ,  $\beta_1\gamma \xRightarrow{*} A\alpha$ , and  $\text{first}(A\alpha) \supseteq \text{first}(A)$  by definition, we have:

$$\text{first}(\beta_1) = \text{first}(\beta_1\gamma) \supseteq \text{first}(A\alpha) \supseteq \text{first}(A) \supseteq \text{first}(w)$$

Because  $\neg\text{nullable}(\beta_2)$  and  $\beta_2\gamma \xRightarrow{*} w$ , we have:

$$\text{first}(\beta_2) = \text{first}(\beta_2\gamma) \supseteq \text{first}(w)$$

Thus

$$\text{first}(\beta_1) \cap \text{first}(\beta_2) \supseteq \text{first}(w) \neq \emptyset$$

This proves that the intersection between the first sets of the RHSs of the two productions  $B \rightarrow \beta_1 \mid \beta_2$  is nonempty, and we have a violation of the LL(1) conditions.

Suppose  $\text{nullable}(\beta_1)$  and  $\neg\text{nullable}(\beta_2)$ . The LL(1) conditions now require  $\text{first}(\beta_1) \cup \text{follow}(B)$  and  $\text{first}(\beta_2)$  to be disjoint. Because  $A \xRightarrow{*} B\gamma$  (assuming no useless symbols or productions),  $\text{follow}(B) \supseteq \text{first}(\gamma)$ . It follows that

$$\text{first}(\beta_1) \cup \text{follow}(B) \supseteq \text{first}(\beta_1) \cup \text{first}(\gamma) = \text{first}(\beta_1\gamma)$$

But then, because  $\beta_1\gamma \xRightarrow{*} A\alpha$  we have

$$\text{first}(\beta_1\gamma) \supseteq \text{first}(A\alpha) \supseteq \text{first}(A) \supseteq \text{first}(w)$$

Thus,  $\text{first}(\beta_1) \cup \text{follow}(B) \supseteq \text{first}(w)$ . But as before, as  $\neg\text{nullable}(\beta_2)$ , we have  $\text{first}(\beta_2) \supseteq \text{first}(w)$ . As  $\text{first}(w) \neq \emptyset$ , we can conclude that these two sets are not disjoint, and therefore the LL(1) conditions are violated.

Suppose instead  $\neg\text{nullable}(\beta_1)$  and  $\text{nullable}(\beta_2)$ . The LL(1) conditions now require  $\text{first}(\beta_1)$  and  $\text{first}(\beta_2) \cup \text{follow}(B)$  to be disjoint. Again,  $\text{follow}(B) \supseteq \text{first}(\gamma)$ . It follows that

$$\text{first}(\beta_2) \cup \text{follow}(B) \supseteq \text{first}(\beta_2) \cup \text{first}(\gamma) = \text{first}(\beta_2\gamma)$$

Then, because  $\beta_2\gamma \xRightarrow{*} w$  we have

$$\text{first}(\beta_2\gamma) \supseteq \text{first}(w)$$

Thus,  $\text{first}(\beta_2) \cup \text{follow}(B) \supseteq \text{first}(w)$ . But given  $\neg\text{nullable}(\beta_1)$ , we have  $\text{first}(\beta_1) \supseteq \text{first}(w)$ . As  $\text{first}(w) \neq \emptyset$ , we can again conclude that these two sets are not disjoint, and therefore the LL(1) conditions are violated.

Thus, no left-recursive grammar can satisfy the LL(1) conditions, which means that a left-recursive grammar first must be transformed into an equivalent grammar that is not left-recursive if we wish to develop an LL(1) parser for the language described by the grammar.

Let us now prove that no ambiguous grammar can be LL(1). Recall that a grammar is ambiguous if a single word  $w$  can be derived in two (or more) essentially different ways, for example there exist two different leftmost derivations for  $w$ .

Assume that a given grammar is ambiguous and pick two different leftmost derivations for some word  $w$  in the language of the grammar. Consider the first place where these derivations differ:

$$S \xRightarrow{*} a_1 \dots a_i A \alpha \Rightarrow a_1 \dots a_i \beta_1 \alpha \xRightarrow{*} a_1 \dots a_i a_{i+1} \dots a_n = w$$

and

$$S \xRightarrow{*} a_1 \dots a_i A \alpha \Rightarrow a_1 \dots a_i \beta_2 \alpha \xRightarrow{*} a_1 \dots a_i a_{i+1} \dots a_n = w$$

Thus there are two productions  $A \rightarrow \beta_1 \mid \beta_2$  in the grammar. Assuming  $a_{i+1} \dots a_n \neq \epsilon$ , we have the following possibilities:

- $a_{i+1} \in \text{first}(\beta_1)$  and  $a_{i+1} \in \text{first}(\beta_2)$ , meaning the LL(1) conditions are violated;
- One of  $\beta_1$  or  $\beta_2$  derives  $\epsilon$ , implying  $a_{i+1} \in \text{follow}(A)$  and  $a_{i+1} \in \text{first}(\beta_1)$  or  $a_{i+1} \in \text{first}(\beta_2)$ , meaning the LL(1) conditions are violated either way;
- Both  $\beta_1$  or  $\beta_2$  derive  $\epsilon$ , a direct violation of the LL(1) conditions.

Assuming  $a_{i+1} \dots a_n = \epsilon$ , then it must be the case that both  $\beta_1$  and  $\beta_2$  are nullable, which violates the LL(1) conditions.

Thus, no ambiguous grammar can satisfy the LL(1) conditions, which means that an ambiguous grammar first must be transformed into an equivalent unambiguous grammar if we wish to develop an LL(1) parser for the described language.

### 10.6.8 Satisfying the LL(1) conditions

Not all grammars satisfy the LL(1) conditions. If we have such a grammar, and we wish to develop an LL(1) parser, the grammar first has to be transformed. In particular, left-recursion and ambiguity must be eliminated because the LL(1) conditions are necessarily violated otherwise (see section 10.6.7). There is no point in computing first and follow sets (for the purpose of constructing a LL(1) parser) before this is done. Of course, transforming a grammar in this way is not always possible: the grammar may be inherently ambiguous, for example (section 8.5). But often transformation into an equivalent grammar satisfying the LL(1) conditions is possible.

Section 8 covered a number of transformations, including ones for eliminating left recursion and disambiguating grammars. Sometimes other transformations are needed. A common problem is the following. Note that the productions are not left recursive, nor would these rules in isolation make a grammar ambiguous as one of the rules derive a word with the terminal  $c$ , and the other one without the terminal  $c$ :

$$S \rightarrow aXbY \mid aXbYcZ$$

This grammar is clearly not suitable for predictive parsing as the first sets for the RHSs of both productions are the same,  $\{a\}$ . But it is also clear that the problem in this case is relatively simple: there is a common prefix of the RHSs of the two productions. Thus, we can try to *postpone* the choice by factoring out the common prefix. That could be enough to satisfy the LL(1) conditions.

Thus, what we need here is *left factoring* (section 8.4). After left factoring:

$$\begin{aligned} S &\rightarrow aXbYS' \\ S' &\rightarrow \epsilon \mid cZ \end{aligned}$$

As it turns out, this is now suitable for LL(1) parsing!

## 10.7 Beyond hand-written parsers: use parser generators

The restriction to LL(1) has a number of disadvantages: In many case a natural grammar like has to be changed to satisfy the LL(1) conditions. This may even be impossible: some context-free *languages* cannot be generated by any LL(1) grammar.

Luckily, there is a more powerful approach called LR(1). LR(1) is a bottom-up method and was briefly discussed in section 10.2. In particular, in contrast to LL(1), LR(1) can handle both left-recursive and right-recursive grammars without modification.

The disadvantage with LR(1) and the related approach LALR(1) (which is slightly less powerful but much more efficient) is that it is very hard to construct LR-parsers by hand. Hence there are automated tools that get the grammar as an input and produce a parser as the output. One of the first of those *parser generators* was YACC for C. Nowadays one can find parser generators for many languages such as JAVA CUP for Java [Hud99] and Happy for Haskell [Mar01].

However, there are also very serious tools based on LL(k) parsing technology, such as ANTLR (ANother Tool for Language Recognition) [Par05], that overcome some of the problems of basic LL parsing and that provides as much automation as any other parser generator tool. It is true that a grammar still may have to be changed a little bit to parseable, but that is true for LR parsing too. In practice, once one become familiar with a tool, it is not that hard to write a grammar in such a way so as to avoid the most common pitfalls from the outset. Additionally, most tools provide mechanics such as declarative specification of disambiguation rules that allows grammars to be written in a natural way without worrying too much about the details of the underlying parsing technology. Today, when it comes to choosing a tool, the underlying parsing technology is probably less important than other factors such as tool quality, supported development languages, feature set, etc.

## 10.8 Exercises

### Exercise 10.1

Consider the following Context-Free Grammar (CFG):

$$\begin{aligned} S &\rightarrow ABB \mid BBC \mid CA \\ A &\rightarrow aA \mid \epsilon \\ B &\rightarrow Bb \mid \epsilon \\ C &\rightarrow cC \mid d \end{aligned}$$

$S$ ,  $A$ ,  $B$ , and  $C$  are nonterminals,  $a$ ,  $b$ ,  $c$ , and  $d$  are terminals, and  $S$  is the start symbol.

1. What is the set  $N_\epsilon$  of *nullable* nonterminals? Provide a brief justification.
2. Systematically compute the *first sets* for all nonterminals, i.e.  $\text{first}(S)$ ,  $\text{first}(A)$ ,  $\text{first}(B)$ , and  $\text{first}(C)$ , by setting up and solving the equations according to the definitions of first sets for nonterminals and strings of grammar symbols. Show your calculations.
3. Set up the subset constraint system that defines the *follow sets* for all nonterminals, i.e.  $\text{follow}(S)$ ,  $\text{follow}(A)$ ,  $\text{follow}(B)$ , and  $\text{follow}(C)$ . Simplify where possible using the law

$$X \subseteq Z \wedge Y \subseteq Z \iff X \cup Y \subseteq Z$$

and the fact that constraints like  $\emptyset \subseteq X$  and  $X \subseteq X$  are trivially satisfied and can be omitted.

4. Solve the subset constraint system for the follow sets from the previous question by finding the smallest sets satisfying the constraints.

### Exercise 10.2

Consider the following Context-Free Grammar (CFG):

$$\begin{aligned} S &\rightarrow AS \mid AB \\ A &\rightarrow aA \mid \epsilon \\ B &\rightarrow BCDb \mid \epsilon \\ C &\rightarrow cD \mid \epsilon \\ D &\rightarrow dC \mid e \end{aligned}$$

$S$ ,  $A$ ,  $B$ ,  $C$ , and  $D$  are nonterminals,  $a$ ,  $b$ ,  $c$ ,  $d$ , and  $e$  are terminals, and  $S$  is the start symbol.

1. What is the set  $N_\epsilon$  of *nullable* nonterminals? Provide a brief justification.
2. Systematically compute the *first sets* for all nonterminals, i.e.,  $\text{first}(S)$ ,  $\text{first}(A)$ ,  $\text{first}(B)$ ,  $\text{first}(C)$ , and  $\text{first}(D)$ , by setting up and solving the equations according to the definitions of first sets for nonterminals and strings of grammar symbols. Show your calculations.

3. Set up the subset constraint system that defines the *follow sets* for all non-terminals; i.e.,  $\text{follow}(S)$ ,  $\text{follow}(A)$ ,  $\text{follow}(B)$ ,  $\text{follow}(C)$ , and  $\text{follow}(D)$ . Simplify where possible using the law

$$X \subseteq Z \wedge Y \subseteq Z \iff X \cup Y \subseteq Z$$

and by removing trivially satisfied constraints such as  $\emptyset \subseteq X$  and  $X \subseteq X$ .

4. Solve the subset constraint system for the follow sets from the previous question by finding the *smallest* sets satisfying the constraints.

## 11 Turing Machines

A *Turing machine* (TM) is a generalization of a PDA that uses a tape instead of a stack. Turing machines are an abstract version of a computer: they have been used to define formally what is *computable*. There are a number of alternative approaches to formalize the concept of computability (for example, the  $\lambda$ -calculus (section 12),  $\mu$ -recursive functions, the von Neumann architecture, and so on) but they all turn out to be equivalent. That this is the case for any reasonable notion of computation is called the *Church-Turing Thesis*.

On the other side there is a generalization of context-free grammars called phrase structure grammars or just grammars. Here we allow several symbols on the left hand side of a production, e.g. we may define the context in which a rule is applicable. Languages definable by grammars correspond precisely to the ones that may be accepted by a Turing machine and those are called *Type-0-languages* or the *recursively enumerable languages* (or *semidecidable languages*)

Turing machines behave differently from the previous machine classes we have seen: they may run forever, without stopping. To say that a language is accepted by a Turing machine means that the TM will stop in an accepting state for each word that is in the language. However, if the word is not in the language the Turing machine may stop in a non-accepting state or loop forever. In this case we can never be sure whether the given word is in the language; i.e., the Turing machine doesn't decide the word problem.

We say that a language is *recursive* (or *decidable*), if there is a TM that accepts it and always stop on any word. There are *type-0-languages* that are not recursive; the most famous one is the *halting problem*. This is the language of encodings of Turing machines that will always stop.

There is no type of grammars that captures all recursive languages (and for theoretical reasons there cannot be one). However there is a subset of recursive languages, called the *context-sensitive languages*, that can be characterized by *context-sensitive grammars*. These are grammars where the left-hand side of a production is always shorter than the right-hand side. Context-sensitive languages on the other hand correspond to linear bounded TMs, that use only a tape whose length can be given by a linear function over the length of the input.

### 11.1 What is a Turing machine?

A Turing machine  $M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$  is:

- A finite set  $Q$  of states;
- A finite set  $\Sigma$  of symbols (the alphabet);
- A finite set  $\Gamma$  of tape symbols s.t.  $\Sigma \subseteq \Gamma$ . This is the case because we use the tape also for the input;
- A transition function

$$\delta \in Q \times \Gamma \rightarrow \{\text{stop}\} \cup Q \times \Gamma \times \{L, R\}$$

The transition function defines how the machine behaves if is in state  $q$  and the symbol on the tape is  $x$ . If  $\delta(q, x) = \text{stop}$  then the machine stops otherwise if  $\delta(q, x) = (q', y, d)$  the machines gets into state  $q'$ , writes  $y$  on the tape (replacing  $x$ ) and moves left if  $d = L$  or right, if  $d = R$ ;

- An initial state  $q_0 \in Q$ ;
- The blank symbol  $B \in \Gamma$  but  $B \notin \Sigma$ . Initially, only a finite section of the tape containing the input is non-blank;
- A set of final states  $F \subseteq Q$ .

In [HMU01] the transition function is defined without an explicit option to stop with type  $\delta \in Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$ . However, they allow  $\delta$  to be undefined which corresponds to our function returning stop.

The above defines deterministic Turing machines; for nondeterministic TMs the type of the transition function is changed to:

$$\delta \in Q \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma \times \{L, R\})$$

In this case, the transition function returning the empty set plays the role of stop. As for finite automata (and unlike for PDAs) there is no difference in the strength of deterministic or nondeterministic TMs.

As for PDAs, we define instantaneous descriptions for Turing machines:  $ID = \Gamma^* \times Q \times \Gamma^*$ . An element  $(\gamma_L, q, \gamma_R) \in ID$  describes a situation where the TM is in state  $Q$ , the non-blank portion of the tape on the left of the head is  $\gamma_L$  and the non-blank portion of the tape on the right, including the square under the head, is  $\gamma_R$ .

We define the next-state relation  $\vdash_M$  similarly to PDAs:

1.  $(\gamma_L, q, x\gamma_R) \vdash_M (\gamma_L y, q', \gamma_R)$  if  $\delta(q, x) = (q', y, R)$
2.  $(\gamma_L z, q, x\gamma_R) \vdash_M (\gamma_L, q', zy\gamma_R)$  if  $\delta(q, x) = (q', y, L)$
3.  $(\epsilon, q, x\gamma_R) \vdash_M (\epsilon, q', By\gamma_R)$  if  $\delta(q, x) = (q', y, L)$
4.  $(\gamma_L, q, \epsilon) \vdash_M (\gamma_L y, q', \epsilon)$  if  $\delta(q, B) = (q', y, R)$
5.  $(\gamma_L z, q, \epsilon) \vdash_M (\gamma_L, q', zy)$  if  $\delta(q, B) = (q', y, L)$
6.  $(\epsilon, q, \epsilon) \vdash_M (\epsilon, q', By)$  if  $\delta(q, B) = (q', y, L)$

The cases 3 to 6 are only needed to deal with the situation of having reached the end of the non-blank part of the tape.

We say that a TM  $M$  *accepts* a word if it goes into an accepting state; i.e., the language of a TM is defined as

$$L(M) = \{w \in \Sigma^* \mid (\epsilon, q_0, w) \vdash_M^* (\gamma_L, q', \gamma_R) \wedge q' \in F\}$$

That is, the TM stops automatically if it goes into an accepting state. However, it may also stop in a non-accepting state if  $\delta$  returns stop. In this case the word is rejected. A TM  $M$  *decides* a language if it accepts it and it never loops (in the negative case).

To illustrate, we define a TM  $M$  that accepts the language

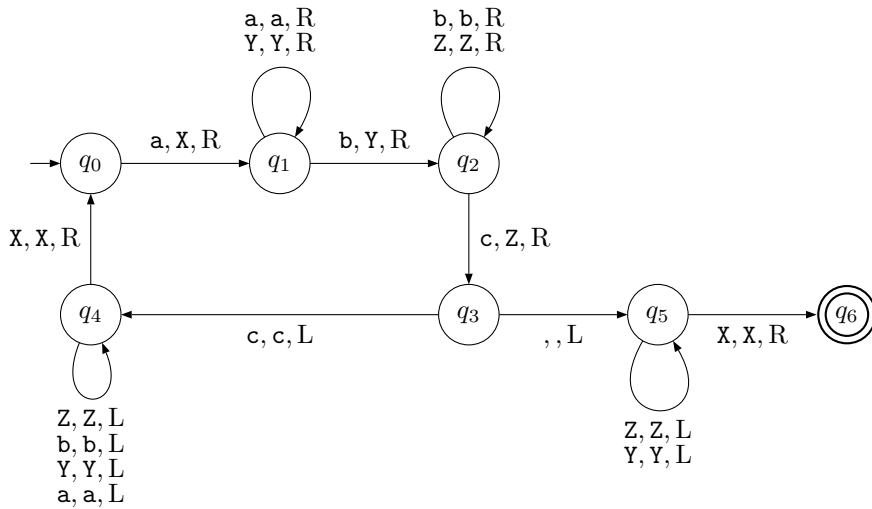
$$L = \{a^n b^n c^n \mid n \in \mathbb{N}\}$$

This is a language that cannot be recognized by a PDA or be defined by a CFG. Let  $M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$  where

$Q$	$= \{q_0, q_1, q_2, q_3, q_4, q_5, q_6\}$	$\delta(q_0, )$	$= (q_6, , R)$
$\Sigma$	$= \{a, b, c\}$	$\delta(q_0, a)$	$= (q_1, X, R)$
$\Gamma$	$= \Sigma \cup \{X, Y, Z, \}$	$\delta(q_1, a)$	$= (q_1, a, R)$
$q_0$	$= q_0$	$\delta(q_1, Y)$	$= (q_1, Y, R)$
$B$	$=$	$\delta(q_1, b)$	$= (q_2, Y, R)$
$F$	$= \{q_6\}$	$\delta(q_2, b)$	$= (q_2, b, R)$
		$\delta(q_2, Z)$	$= (q_2, Z, R)$
		$\delta(q_2, c)$	$= (q_3, Z, R)$
		$\delta(q_3, )$	$= (q_5, , L)$
		$\delta(q_3, c)$	$= (q_4, c, L)$
		$\delta(q_4, Z)$	$= (q_4, Z, L)$
		$\delta(q_4, b)$	$= (q_4, b, L)$
		$\delta(q_4, Y)$	$= (q_4, Y, L)$
		$\delta(q_4, a)$	$= (q_4, a, L)$
		$\delta(q_4, X)$	$= (q_0, X, R)$
		$\delta(q_5, Z)$	$= (q_5, Z, L)$
		$\delta(q_5, Y)$	$= (q_5, Y, L)$
		$\delta(q_5, X)$	$= (q_6, X, R)$
		$\delta(q, x)$	$= \text{stop} \quad \text{everywhere else}$

The machine replaces an  $a$  by  $X$  ( $q_0$ ) and then looks for the first  $b$  replaces it by  $Y$  ( $q_1$ ) and looks for the first  $c$  and replaces it by a  $Z$  ( $q_2$ ). If there are more  $c$ 's left it moves left to the next  $a$  ( $q_4$ ) and repeats the cycle. Otherwise it checks whether there are no  $a$ 's and  $b$ 's left ( $q_5$ ) and if so goes into an accepting state ( $q_6$ ).

Graphically the machine can be represented by the following transition diagram, where the edges are labelled by (read-symbol, write-symbol, move-direction):





E.g. consider the sequence of IDs on **aabbcc**:

$$\begin{aligned}
(\epsilon, q_0, \mathbf{aabbcc}) &\vdash (\mathbf{X}, q_1, \mathbf{abbcc}) \\
&\vdash (\mathbf{Xa}, q_1, \mathbf{bbcc}) \\
&\vdash (\mathbf{XaY}, q_2, \mathbf{bcc}) \\
&\vdash (\mathbf{XaYb}, q_2, \mathbf{cc}) \\
&\vdash (\mathbf{XaYbZ}, q_3, \mathbf{c}) \\
&\vdash (\mathbf{XaYb}, q_4, \mathbf{Zc}) \\
&\vdash (\mathbf{XaY}, q_4, \mathbf{bZc}) \\
&\vdash (\mathbf{Xa}, q_4, \mathbf{YbZc}) \\
&\vdash (\mathbf{X}, q_4, \mathbf{aYbZc}) \\
&\vdash (\epsilon, q_4, \mathbf{XaYbZc}) \\
&\vdash (\mathbf{X}, q_0, \mathbf{aYbZc}) \\
&\vdash (\mathbf{XX}, q_1, \mathbf{YbZc}) \\
&\vdash (\mathbf{XXY}, q_1, \mathbf{bZc}) \\
&\vdash (\mathbf{XXYY}, q_2, \mathbf{Zc}) \\
&\vdash (\mathbf{XXYYZ}, q_2, \mathbf{c}) \\
&\vdash (\mathbf{XXYYZZ}, q_3, \epsilon) \\
&\vdash (\mathbf{XXYYZ}, q_5, \mathbf{Z}) \\
&\vdash (\mathbf{XXYY}, q_5, \mathbf{ZZ}) \\
&\vdash (\mathbf{XXY}, q_5, \mathbf{YZZ}) \\
&\vdash (\mathbf{XX}, q_5, \mathbf{YYZZ}) \\
&\vdash (\mathbf{X}, q_5, \mathbf{XYYZZ}) \\
&\vdash (\mathbf{XX}, q_6, \mathbf{YYZZ})
\end{aligned}$$

We see that  $M$  accepts **aabbcc**. Because  $M$  never loops, it does actually decide  $L$ .

## 11.2 Grammars and context-sensitivity

Let us define grammars  $G = (N, T, P, S)$  in the same way as context-free grammars were defined before, with the *only* difference that there may now be *several* symbols on the left-hand side of a production; i.e.,  $P \subseteq (N \cup T)^+ \times (N \cup T)^*$ . Here  $(N \cup T)^+$  means that at least one symbol has to present. The relation derives  $\Rightarrow_G$  (and  $\xRightarrow{*}_G$ ) is defined as before:

$$\begin{aligned}
&\xRightarrow{*}_G \subseteq (N \cup T)^* \times (N \cup T)^* \\
\alpha\beta\gamma \xRightarrow{*}_G \alpha\beta'\gamma &\iff \beta \rightarrow \beta' \in P
\end{aligned}$$

Also as before, the language of  $G$  is defined as:

$$L(G) = \{w \in T^* \mid S \xRightarrow{*}_G w\}$$

We say that a grammar is context-sensitive (or type 1) if the right-hand side of a production is at least as long as the left-hand side. That is, for each  $\alpha \rightarrow \beta \in P$ , we have  $|\beta| \geq |\alpha|$ .

Here is an example of a context-sensitive grammar  $G = (N, T, P, S)$  with  $L(G) = \{a^n b^n c^n \mid n \in \mathbb{N} \wedge n \geq 1\}$  where

- $N = \{S, B, C\}$
- $T = \{a, b, c\}$
- $P$  is the set of productions:

$$\begin{array}{ll} S & \rightarrow aSBC \\ S & \rightarrow aBC \\ aB & \rightarrow ab \\ CB & \rightarrow BC \\ bB & \rightarrow bb \\ bC & \rightarrow bc \\ cC & \rightarrow cc \end{array}$$

- $S$  is the start symbol

We present without proof:

**Theorem 11.1** *For a language  $L \subseteq T^*$  the following is equivalent:*

1.  $L$  is accepted by a Turing machine  $M$ ; i.e.,  $L = L(M)$
2.  $L$  is given by a grammar  $G$ ; i.e.,  $L = L(G)$

**Theorem 11.2** *For a language  $L \subseteq T^*$  the following is equivalent:*

1.  $L$  is accepted by a Turing machine  $M$ ; i.e.,  $L = L(M)$  such that the length of the tape is bounded by a linear function in the length of the input; i.e.,  $|\gamma_L| + |\gamma_R| \leq f(x)$  where  $f(x) = ax + b$  with  $a, b \in \mathbb{N}$ .
2.  $L$  is given by a context-sensitive grammar  $G$ ; i.e.,  $L = L(G)$

### 11.3 The halting problem

Turing showed that there are languages that are accepted by a TM (i.e., type 0 languages) but that are undecidable. The technical details of this construction are involved but the basic idea is simple and is closely related to Russell's paradox, which we have seen in MCS.

Let's fix a simple alphabet  $\Sigma = \{0, 1\}$ . As computer scientist we are well aware that everything can be coded up in bits and hence we accept that there is an encoding of TMs in binary. That is, we can agree on a way to express every TM  $M$  as a string of bits  $\lceil M \rceil \in \{0, 1\}^*$ . We assume that the string contains its own length at the beginning, so that we know when the encoding ends. This allows us to put both the TM and an input for it, one after the other, on the same tape. We can determine when the encoding of the machine ends and the subsequent input on starts.

Now we define the following language

$$L_{\text{halt}} = \{\lceil M \rceil w \mid M \text{ halts on input } w.\}$$

It is easy (although the details are quite daunting) to define a TM that accepts this language: we just simulate  $M$  and accept if  $M$  stops. However, Turing showed that there is no TM that decides this language.

Let us prove this by assuming that the language is decidable and deriving a contradiction from it. Suppose there is a TM  $H$  that decides  $L$ . That is, when run on a word  $v$ ,  $H$  always terminates and the final state is accepting if and only if  $v$  has the form  $v = \lceil M \rceil w$  for some machine  $M$  and input  $w$  and when  $M$  is run on  $w$  it terminates. In all other cases, if  $v$  is in such form but  $M$  does not terminate on  $w$  or if  $v$  is not in that form at all,  $H$  will terminate in a rejecting state.

Now using  $H$  we construct a new TM  $F$  that is a bit obnoxious. When we run  $F$  on input  $x$ , it computes  $H$  on the duplicated input  $xx$ . If  $H$  says *yes* (it accepts  $xx$ ), then  $F$  goes into an infinite loop; otherwise, if  $H$  says *no* (it rejects  $xx$ ),  $F$  stops.

What happens if we run  $F$  on its own code  $\lceil F \rceil$ ? Will it terminate or loop forever? If we consider each of the two possibilities, we get the contradictory conclusion that the opposite should be true: if we assume that it terminates, we can prove that it must loop; if we assume that it loops, we can prove that it must terminate!

Let us assume  $F$  on  $\lceil F \rceil$  terminates. By definition of  $F$ , this happens only if  $H$  applied to  $\lceil F \rceil \lceil F \rceil$  says *no*. But by definition of  $H$ , this means that  $F$  on  $\lceil F \rceil$  loops, contradicting the assumption.

Let us then assume that  $F$  on input  $\lceil F \rceil$  loops. This happens only if  $H$  applied to  $\lceil F \rceil \lceil F \rceil$  says *yes*. But this means that  $F$  on  $\lceil F \rceil$  terminates, again contradicting the assumption.

We reach a contradiction on both possible behaviours of  $F$ : the machine  $F$  cannot exist. We must conclude that our assumption that there is a TM  $H$  that decides  $L_{\text{halt}}$  is false. We say  $L_{\text{halt}}$  is undecidable.

We have shown that there is no Turing machine that can decide whether other Turing machines halt. Is this a specific shortcoming of TMs or a universal limitation of all computing models? Maybe we could find a more powerful programming language that overcomes this problem? It turns out that all computational formalisms (i.e., programming languages) that have actually been implemented are equal in power and can be simulated by each other. Beside TMs, all modern programming languages, the  $\lambda$ -calculus,  $\mu$ -recursive functions and many others were proved to be equivalent.

The statement that all models of computability are equivalent is called the *Church-Turing thesis* because it was first formulated by Alonzo Church and Alan Turing in the 1930s. This is discussed further in section 12.6.

## 11.4 Recursive and recursively enumerable sets

When we work with Turing Machines, there is a distinction between a language being accepted or decided.

A machine  $M$  *accepts* a language  $L$  if, whenever we run  $M$  on an input  $w$ , the computation terminates in an accepting state if and only if  $w \in L$ . We say that  $M$  *decides*  $L$  if it accepts it and always terminates.

The difference is that a machine that just accepts a language but doesn't decide it may run forever on some words that do not belong to the language.

If that happens, we have to wait forever to discover whether the word is in the language or not.

**Definition 11.3** A language  $L$  is recursively enumerable if it is accepted by a Turing Machine.

The terminology comes from a different characterization of such languages. We can define a *recursively enumerable set* as a collection of values that can be produced by a total algorithm. Whenever we have a computable function  $f : \mathbb{N} \rightarrow A$  from the natural numbers to some set  $A$ , we say that the range of the function,  $U = \{a \in A \mid \exists n : \mathbb{N}, f(n) = a\}$  is recursively enumerable. We can also write  $U = \{f(n) \mid n \in \mathbb{N}\}$ . The idea is that we can enumerate all elements of  $U$  by computing  $f(0)$ ,  $f(1)$ ,  $f(2)$ , etc. In the case of languages, the set  $A$  is  $\Sigma^*$ . We can prove that being a recursively enumerable subset of  $\Sigma^*$  is equivalent to be accepted by a Turing Machine.

**Definition 11.4** A language  $L$  is recursive if there is a Turing machine that accepts it and always terminate.

In general, a *recursive set* is a subset of some set  $A$  for which we can effectively determine membership. Whenever we have a computable function  $g : A \rightarrow \text{Bool}$  from some set  $A$  to Booleans, we say that the preimage of **true**,  $\{a \in A \mid g(a) = \text{true}\}$  is recursive.

**Theorem 11.5** A subset  $U \subseteq A$  is recursive if and only if both  $U$  and its complement  $\bar{U}$  are recursively enumerable.

**Proof.** We won't see a formal proof using Turing Machines, but I'll outline the idea.

- In one direction, assume that  $U$  is recursive. So we have a function  $g : A \rightarrow \text{Bool}$  that decides it.

We can construct  $f : \mathbb{N} \rightarrow A$  by using some enumeration of all of  $A$  (we assume it is computably countable, otherwise it won't make sense to talk about algorithms on it). When computing the values of  $f$  we keep an index  $n$ ; we go through all the elements of  $A$  one by one, whenever we find an  $a$  such that  $g(a) = \text{true}$ , we set  $f(n) = a$  and we increase  $n$  by one, now looking for the next element of  $A$  satisfying  $g$ . In this way we construct an  $f$  that generates all elements of  $U$ , so we proved that  $U$  is recursively enumerable.

In a similar way we can show that  $\bar{U}$  is recursively enumerable by systematically searching for the elements  $a \in A$  for which  $g(a) = \text{false}$ .

- In the other direction, assume that both  $U$  and  $\bar{U}$  are recursively enumerable. We must prove that  $U$  is recursive.

Saying that  $U$  is recursively enumerable means that there is a computable function  $f_1 : \mathbb{N} \rightarrow A$  such that  $U = \{f_1(a) \mid n : \mathbb{N}\}$ . Saying that  $\bar{U}$  is also recursively enumerable means that there is a computable function  $f_2 : \mathbb{N} \rightarrow A$  such that  $\bar{U} = \{f_2(a) \mid n : \mathbb{N}\}$ .

From  $f_1$  and  $f_2$  we can construct a function  $g : A \rightarrow \text{Bool}$  that decides  $U$ . For a given element  $a \in A$ , we can search whether it belongs to  $U$  by running the two functions  $f_1$  and  $f_2$  repeatedly in parallel:

- run  $f_1$  on 0, if  $f_1(0) = a$  then we know that  $a \in U$  and we terminate giving the answer  $g(a) = \text{true}$ ;
- run  $f_2$  on 0, if  $f_2(0) = a$  then we know that  $a \in \overline{U}$  and we terminate giving the answer  $g(a) = \text{false}$ ;
- run  $f_1$  on 1, if  $f_1(1) = a$  then we terminate with answer  $g(a) = \text{true}$ ;
- run  $f_2$  on 1, if  $f_2(1) = a$  then we terminate with answer  $g(a) = \text{false}$ ;
- run  $f_1$  on 2, if  $f_1(2) = a$  then we terminate with answer  $g(a) = \text{true}$ ;
- run  $f_2$  on 2, if  $f_2(2) = a$  then we terminate with answer  $g(a) = \text{false}$ ;
- continue until you find an  $n$  such that either  $f_1(n) = a$  or  $f_2(n) = a$ .

Because every  $a$  belongs to either  $U$  or  $\overline{U}$ , we know for sure that this process will always terminate and produce the correct answer.

We have constructed a function  $g$  that decides  $U$ , so  $U$  is recursive.

□

In the previous section we proved that the Halting Problem is undecidable, therefore the set  $L_{\text{halt}} = \{[M]w \mid M \text{ halts on input } w\}$  is not recursive.

It is easy to see that it is recursively enumerable: we just have to run  $M$  and see if it terminates.

There are many other problems that are semi-decidable but not decidable, that is, they can be expressed by a recursively enumerable language that is not recursive. We can prove this by reducing the Halting Problem (or any other already known undecidable problem) to them.

**Definition 11.6** *A language  $L_1$  is reducible to a language  $L_2$  if there is a Turing Machine  $M$  such that:*

- $M$  terminates on every input;
- If we run  $M$  on an input  $w$ , the computation terminates with a word  $v$  on the tape such that  $w \in L_1$  if and only if  $v \in L_2$ .

So  $L_1$  is reducible to  $L_2$  if we can translate (by a computable function) every question of membership of  $L_1$  to an equivalent membership question of  $L_2$ .

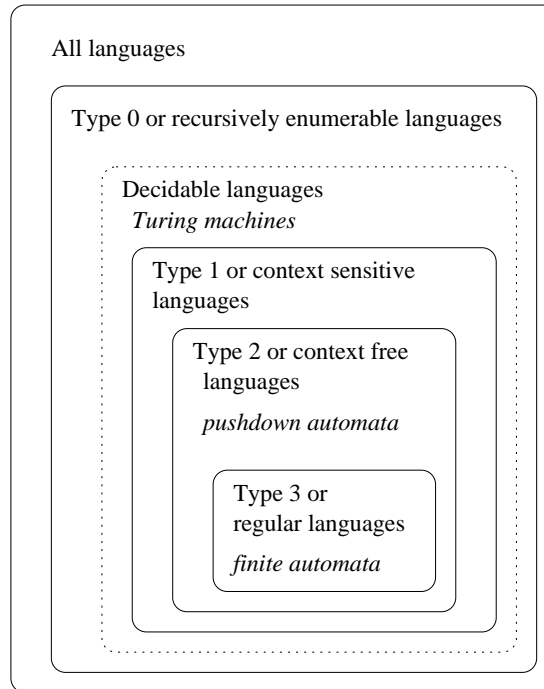
**Theorem 11.7** *If  $L_{\text{halt}}$  is reducible to some language  $L$ , then  $L$  is not decidable.*

**Proof.** Suppose, towards a contradiction, that  $L$  is recursive. Then we have a Turing Machine  $M$  that decides  $L$ . We also know that there is a Turing Machine  $M'$  that translates instances of the Halting Problem to  $L$ . But then, by composing  $M'$  and  $M$ , we would be able to construct a machine  $H$  that decides the halting problem. But we know this is impossible. Therefore our assumption that  $L$  is recursive must be false. □

In general, if we know that a set  $V$  is not recursive and have another set  $U$  that we suspect is also not recursive, we can prove this fact by reducing  $V$  to  $U$ . Be very careful in the direction of the reduction: We reduce the problem that we already know to be undecidable to the one for which we want to prove undecidability.

## 11.5 Back to Chomsky

At the end of the course we should have another look at the Chomsky hierarchy, which classifies languages based on subclasses of grammars, or equivalently by different types of automata that recognize them



We have worked our way from the bottom to the top of the hierarchy: starting with finite automata, computation with fixed amount of memory via pushdown automata (finite automata with a stack), to Turing machines (finite automata with a tape). Correspondingly we have introduced different grammatical formalisms: regular expressions, context-free grammars and grammars.

Note that at each level there are languages that are on the next level but not on the previous:  $\{a^n b^n \mid n \in \mathbb{N}\}$  is Type 2 but not Type 3,  $\{a^n b^n c^n\}$  is Type 1 but not Type 2, and the Halting problem is Type 0 but not Type 1.

We could have gone the other way: starting with Turing machines and grammars and then introducing restrictions: Turing machines that only use their tapes as a stack describe Type 2 languages; Turing machines that never use the tape apart for reading the input describe Type 3 languages. Similarly, we have seen that context-free grammars restricted in specific ways describe precisely the regular languages (section 7.3).

Chomsky introduced his hierarchy as a classification of grammars; the relation to automata was only observed a bit later. This may be the reason why he introduced the Type-1 level, which is not so interesting from an automata point of view (unless you are into computational complexity; i.e., resource use, here linear use of memory). It is also the reason why on the other hand the decidable languages do not constitute a level: there is no corresponding grammatical formalism (we can even prove this).

## 11.6 Exercises

### Exercise 11.1

Consider the Turing Machine defined formally by:

$$M = (Q, \Sigma, \Gamma, \delta, q_0, \sqcup, F) \quad \text{where} \quad \begin{aligned} Q &= \{q_0, q_1, q_2, q_3, q_4, q_5\} \\ F &= \{q_5\} \\ \Sigma &= \{a, b\} \\ \Gamma &= \{a, b, X, Y, \sqcup\} \end{aligned}$$

with the transition function defined by:

$$\begin{aligned} \delta(q_0, a) &= (q_3, X, L) \\ \delta(q_0, x) &= (q_0, x, R) \quad \text{for } x \in \{b, X, Y\} \\ \delta(q_0, \sqcup) &= (q_2, \sqcup, L) \\ \delta(q_1, b) &= (q_4, Y, L) \\ \delta(q_1, x) &= (q_1, x, R) \quad \text{for } x \in \{a, X, Y\} \\ \delta(q_2, \sqcup) &= (q_5, \sqcup, R) \\ \delta(q_2, x) &= (q_2, x, L) \quad \text{for } x \in \{X, Y\} \\ \delta(q_3, \sqcup) &= (q_1, \sqcup, R) \\ \delta(q_3, x) &= (q_3, x, L) \quad \text{for } x \in \{a, b, X, Y\} \\ \delta(q_4, \sqcup) &= (q_0, \sqcup, R) \\ \delta(q_4, x) &= (q_4, x, L) \quad \text{for } x \in \{a, b, X, Y\} \end{aligned}$$

1. Draw  $M$  graphically as a transition diagram.
2. Write down the sequence of instantaneous descriptions when starting the machine  $M$  with input  $baab$ . Is this word accepted or rejected?
3. Write down the sequence of instantaneous descriptions when starting the machine  $M$  with input  $aba$ . Is this word accepted or rejected?
4. What is the language accepted by  $M$ ?

### Exercise 11.2

Construct a Turing Machine that, when run on any word in  $\{a, b\}^*$ , rearranges the symbols so all the  $a$ s come before all the  $b$ s. For example:

$$(\epsilon, q_0, aababba) \vdash^* (aaaabbb, q_f, \epsilon) \quad (\epsilon, q_0, babbbba) \vdash^* (aaabbbb, q_f, \epsilon)$$

[Hint: Your machine could look for pairs of symbols in the order ‘ $ba$ ’ and swap them to ‘ $ab$ ’; keep repeating this until there are no such pairs left.]

1. Draw the machine as a transition diagram.
2. Give a formal definition of the machine, defining the transition function.

### Exercise 11.3

Suppose that you have the following information about ten formal problems (expressed as languages)  $P_1, P_2, \dots, P_{10}$ :

- The Halting Problem is reducible to  $P_1$ ;
- $P_2$  is reducible to  $P_1$ ;
- $P_1$  is reducible to  $P_5$ ;
- The complement of  $P_3$  is recursively enumerable;
- $P_4$  is recursively enumerable;
- $P_4$  is reducible to the normalization problem for  $\lambda$ -calculus;
- $P_4$  is reducible to  $P_5$ ;
- $P_6$  is not recursively enumerable;
- $P_7$  is recursive;
- $P_7$  is reducible to  $P_6$ ;
- $P_7$  is reducible to  $P_8$ ;
- The complement of  $P_8$  is not recursively enumerable;
- The normalization problem for  $\lambda$ -calculus is reducible to  $P_9$ ;
- $P_{10}$  is reducible to  $P_3$ ;
- $P_{10}$  is reducible to  $P_4$ .

Note that the normalization problem for  $\lambda$ -calculus (see section 12) is *undecidable*. Given this information, which of the languages  $P_1$  to  $P_{10}$  are:

1. Undecidable?
2. Recursively enumerable?
3. Recursive?

Note that the same language may be in more than one of these three classes. Some others may be in none. Justify your answers by explaining how your conclusions follow from the given information.



## 12 $\lambda$ -Calculus

In traditional imperative programming (like C, Java, Python), we have a clear distinction between *programs*, that are sequences of instructions that are executed sequentially, and *data*, that are values given in input, stored in memory, manipulated during computation and returned as output. Programs and data are distinct and are kept separated. Programs are not modified during computation. (It is in theory possible to do it, because programs are stored in memory like any other data. However, it is difficult and dangerous.)

*Functional Programming* is a different paradigm of computation in which there is no distinction between programs and data. Both are represented by terms/expressions belonging to the same language. Computation consists in the reduction of terms to normal form. That includes terms that represent functions, that is, the programs themselves.

The pure realization of this idea is the  $\lambda$ -calculus. It is a pure theory of functions with only one kind of objects:  $\lambda$ -terms. They represent both data structures and programs.

The main idea is the definition of functions by *abstraction*. For example, we may define a function  $f$  on numbers by saying that  $f(x) = x^2 + 3$ . By this we mean that any argument to the function, represented by the variable  $x$ , is squared and added to 3. The use of variables is different from imperative programming:  $x$  is just a place-holder to denote any possible value, while in imperative programming variables represent memory locations containing values that can be modified.

We can specify the function  $f$  alternatively with the *mapping* notation:

$$x \mapsto x^2 + 3.$$

This is written in  $\lambda$ -notation as:  $f = \lambda x. x^2 + 3$ . (In the functional programming language Haskell, it is `\x -> x^2+3`.)

While abstraction is the operation to define a new function, computing it on a specific argument is called *application*. We indicate it simply by juxtaposition:

$$f\ 5 = (\lambda x. x^2 + 3)\ 5 \rightsquigarrow 5^2 + 3 \rightsquigarrow^* 28.$$

As the example shows, the application of a  $\lambda$ -abstraction to an argument is computed by replacing the abstraction variable with the argument. This is called  $\beta$ -reduction and it is the basic computation step of  $\lambda$ -calculus.

The  $\lambda$ -notation is convenient to define functions, but you may think that the actual computation work is done by the operations used in the body of the abstraction: squaring and adding 3. However, the  $\lambda$ -calculus is *a theory of pure functions*: terms are constructed using only abstraction and application, there are no other basic operations. At first, this looks like a rather useless system: no numbers, no arithmetic operations, no data structures, no programming primitives. The surprising fact is that we don't really need them. We don't need numbers (5, 3) and we don't need operations ( $-^2$ ,  $+$ ). They all can be defined as purely functional constructions, built using only abstraction and application!

### 12.1 Syntax of $\lambda$ -calculus

The language of  $\lambda$ -calculus is extremely simple, we start with variables and construct terms using only abstraction and computation. It's BNF definition is

as follows (assume  $x, y, z$  range over a given infinite set of variable names).

$$\begin{array}{ll} t := & x \mid y \mid z \mid \dots & \text{variable names} \\ & \mid \lambda x.t & \text{abstraction} \\ & \mid tt & \text{application.} \end{array}$$

The  $\beta$ -reduction relation on term is defined, for every pair of terms  $t_1$  and  $t_2$  as:

$$(\lambda x.t_1) t_2 \rightsquigarrow_{\beta} t_1[x := t_2].$$

The left-hand side means: in  $t_1$  substitute all occurrences of variable  $x$  with the term  $t_2$ . Substitution is actually quite tricky and its precise definition is a bit more complex than replacing every occurrences of  $x$  with  $t_2$ . One has to be careful to manage variable occurrences properly.

We need some intermediate concept. The first is  $\alpha$ -equivalence and it says that, because the variable in an abstraction is just a place holder for an argument, the names of abstracted variables does not matter. For example, the simplest function we can define is the identity  $\text{id} := \lambda x.x$  which takes an argument  $x$  and returns it unchanged. Clearly, if we use a different variable name,  $\lambda y.y$ , we get exactly the same function. We say that the two terms (and any using different variables) are  $\alpha$ -equivalent:

$$\lambda x.x =_{\alpha} \lambda y.y =_{\alpha} \lambda z.z =_{\alpha} \dots$$

We are free to change the name of the abstracted variable any way we like. However, we have to be careful to avoid *variable capture*. If the body of the abstraction contains other variables than the abstracted one, we can't change the name to those:

$$\lambda x.y(xz) =_{\alpha} \lambda w.y(wz) \neq_{\alpha} \lambda y.y(yz) \neq_{\alpha} \lambda z.y(zz).$$

The reason for this restriction is that changing the name of the abstracted variable from  $x$  to either  $y$  or  $z$  in this example would capture the occurrence of that variable which was *free* in the original term (not bound by a  $\lambda$ -abstraction).

Formally, we define set  $\text{FV}(t)$  of the variables that occur free in the term  $t$ , by recursion on the structure of  $t$ :

$$\begin{aligned} \text{FV}(x) &= \{x\} \\ \text{FV}(\lambda x.t) &= \text{FV}(t) \setminus \{x\} \\ \text{FV}(t_1 t_2) &= \text{FV}(t_1) \cup \text{FV}(t_2). \end{aligned}$$

Another way in which a variable can be incorrectly captured is when we perform a substitution that puts a term under an abstraction that may bind some of its variables:

$$(\lambda x.\lambda y.xy)(yz) \rightsquigarrow_{\beta} (\lambda y.xy)[x := (yz)] \neq \lambda y.(yz)y.$$

If we replace  $x$  with  $(yz)$  in this way, the occurrence of the variable  $y$  (which was free before performing the  $\beta$ -reduction) become bound. This is incorrect. We should rename the abstraction variable before performing the substitution:

$$(\lambda y.xy)[x := (yz)] =_{\alpha} (\lambda w.xw)[x := (yz)] = \lambda w.(yz)w.$$

To avoid problems with variable capture, we adopt the *Barendregt variable convention*: before performing substitution (or any other operation on terms), change the names of the abstracted variables so they are different from the free variables and from each other.

With this convention, we can give a precise definition of substitution by recursion on the structure of terms:

$$\begin{aligned} x[x := t_2] &= t_2 \\ y[x := t_2] &= y \quad \text{if } y \neq x \\ (\lambda y. t_1)[x := t_2] &= \lambda y. t_1[x := t_2] \\ (t_0 t_1)[x := t_2] &= t_0[x := t_2] t_1[x := t_2]. \end{aligned}$$

In the third case, the variable convention ensures that the variable  $y$  and the bound variables in  $t_2$  have already been renamed so they avoid captures. In more traditional formulations, one would add the requirements: “provided that  $y \neq x$  and  $y$  doesn’t occur free in  $t_2$ ”.

A part from some complication about substitution, the  $\lambda$ -calculus is extremely simple. It seems at first surprising that we can actually do any serious computation with it at all. But it turns out that all computable functions can be represented by  $\lambda$ -terms. We see some simple function in this section and we will discover how to represent data structures in the next.

For convenience, we use some conventions that allow us to save on parentheses.

- $\lambda$ -abstraction associates to the right, so we write  $\lambda x. \lambda y. x$  for  $\lambda x. (\lambda y. x)$ ;
- Application associates to the left, so we write  $(t_1 t_2 t_3)$  for  $((t_1 t_2) t_3)$ ;
- We can use a single  $\lambda$  symbol followed by several variables to mean consecutive abstractions, so we write  $\lambda x y. x$  for  $\lambda x. \lambda y. x$ .

Here are three very simple functions implemented as  $\lambda$ -terms:

- The identity function  $\text{id} := \lambda x. x$ . When applied to an argument it simply returns it unchanged:

$$\text{id } t = (\lambda x. x) t \rightsquigarrow x[x := t] = t.$$

- The first projection function  $\lambda x. \lambda y. x$ . When applied to two arguments, it returns the first:

$$(\lambda x. \lambda y. x) t_1 t_2 \rightsquigarrow (\lambda y. x)[x := t_1] t_2 = (\lambda y. t_1) t_2 \rightsquigarrow t_1[y := t_2] = t_1.$$

Remember, in reading this reduction sequence, that we are adopting the variable convention, so the variable  $y$  doesn’t occur free in  $t_1$ .

- The second projection function  $\lambda x. \lambda y. y$ . When applied to two arguments, it returns the second:

$$(\lambda x. \lambda y. y) t_1 t_2 \rightsquigarrow (\lambda y. y)[x := t_1] t_2 = (\lambda y. y) t_2 \rightsquigarrow y[y := t_2] = t_2.$$

We can also say that the second projection is the function that, when applied to an argument  $t_1$ , returns the identity function  $\lambda y. y$ .

## 12.2 Church numerals

So far we have seen only some very basic functions that only return some of their arguments unchanged. How can we define more interesting computations? And first of all, how can we represent values and data structures? It is in fact possible to represent any kind of data by some  $\lambda$ -term.

Let's start by representing natural numbers. Their encodings in  $\lambda$ -calculus are called *Church Numerals*:

$$\begin{aligned}\bar{0} &:= \lambda f.\lambda x.x \\ \bar{1} &:= \lambda f.\lambda x.f\ x \\ \bar{2} &:= \lambda f.\lambda x.f\ (f\ x) \\ \bar{3} &:= \lambda f.\lambda x.f\ (f\ (f\ x)) \\ &\dots\end{aligned}$$

A numeral  $\bar{n}$  is a function that takes two arguments, denoted by the variables  $f$  and  $x$ , and applies  $f$  sequentially  $n$  times to  $x$ .

What is important is that we assign to every number a distinct  $\lambda$ -term in a uniform way. We must choose our representation so it is easy to represent arithmetic operations. The idea of Church numerals is nicely conceptual: numbers are objects that we use to count things, so we can define them as the counters of repeated application of a function.

Let's see if this representation is convenient from the programming point of view: can we define basic operations on it?

Let's start with the successor function, that increases a number by one:

$$\text{succ} := \lambda n.\lambda f.\lambda x.f\ (n\ f\ x).$$

Let's test if it works on an example: if we apply it to  $\bar{2}$  we should get  $\bar{3}$ :

$$\begin{aligned}\text{succ}\ \bar{2} &= (\lambda n.\lambda f.\lambda x.f\ (n\ f\ x))\ \bar{2} \\ &\rightsquigarrow \lambda f.\lambda x.f\ (\bar{2}\ f\ x) = \lambda f.\lambda x.f\ ((\lambda f.\lambda x.f\ (f\ x))\ f\ x) \\ &\rightsquigarrow \lambda f.\lambda x.f\ ((\lambda x.f\ (f\ x))[f := f]\ x) = \lambda f.\lambda x.f\ ((\lambda x.f\ (f\ x))\ x) \\ &\rightsquigarrow \lambda f.\lambda x.f\ ((f\ (f\ x))[x := x]) = \lambda f.\lambda x.f\ (f\ (f\ x)) = \bar{3}.\end{aligned}$$

We have explicitly marked the substitutions in this reduction sequence: they are both trivial, substituting  $f$  with itself and  $x$  with itself. From now on, we'll do the substitutions on the fly, without marking them.

Other arithmetic operations can be defined by simple terms:

$$\begin{aligned}\text{plus} &:= \lambda m.\lambda n.\lambda f.\lambda x.m\ f\ (n\ f\ x) \\ \text{mult} &:= \lambda m.\lambda n.\lambda f.m\ (n\ f) \\ \text{exp} &:= \lambda m.\lambda n.n\ m.\end{aligned}$$

Verify by yourself that these terms correctly implement the addition, multiplication and exponentiation functions. For example:

$$\begin{aligned}\text{plus}\ \bar{2}\ \bar{3} &\rightsquigarrow^* \bar{5} \\ \text{mult}\ \bar{2}\ \bar{3} &\rightsquigarrow^* \bar{6} \\ \text{exp}\ \bar{2}\ \bar{3} &\rightsquigarrow^* \bar{8} \\ \text{exp}\ \bar{2}\ \bar{3} &\rightsquigarrow^* \bar{9}.\end{aligned}$$

**Exercise.** Define a term `isZero` that tests if a Church numeral is  $\bar{0}$  or a successor. It should have the following reduction behaviour:

$$\text{isZero } \bar{0} \rightsquigarrow^* \text{true}; \quad \text{isZero } (\text{succ } t) \rightsquigarrow^* \text{false} \quad \text{for every term } t.$$

Surprisingly, two other basic functions are much more difficult to define: the predecessor and the (cut-off) subtraction functions. Try to define two terms `pred` and `minus` such that:

$$\begin{aligned} \text{pred } (\text{succ } \bar{n}) &\rightsquigarrow^* \bar{n} \\ \text{pred } \bar{0} &\rightsquigarrow^* \bar{0} \\ \text{minus } \bar{n} \bar{m} &\rightsquigarrow^* \overline{n - m} \quad \text{if } n \geq m \\ \text{minus } \bar{n} \bar{m} &\rightsquigarrow^* \bar{0} \quad \text{if } n < m. \end{aligned}$$

### 12.3 Other data structures

Other data types can be encoded in the  $\lambda$ -calculus.

**Booleans** For truth values we may choose the first and second projects that we defined earlier:

$$\begin{aligned} \text{true} &:= \lambda x. \lambda y. x \\ \text{false} &:= \lambda x. \lambda y. y. \end{aligned}$$

We must show how to compute the logical operators. For example, conjunction can be defined as follows:

$$\text{and} := \lambda a. \lambda b. a \ b \ \text{false}.$$

Let's verify that it give the correct results when applied to Boolean values:

$$\begin{aligned} \text{and true true} &= (\lambda a. \lambda b. a \ b \ \text{false}) \ \text{true} \ \text{true} \\ &\rightsquigarrow^* \text{true true false} = (\lambda x. \lambda y. x) \ \text{true} \ \text{false} \rightsquigarrow^* \text{true} \\ \text{and true false} &= (\lambda a. \lambda b. a \ b \ \text{false}) \ \text{true} \ \text{false} \\ &\rightsquigarrow^* \text{true false false} = (\lambda x. \lambda y. x) \ \text{false} \ \text{false} \rightsquigarrow^* \text{false} \\ \text{and false } t &= (\lambda a. \lambda b. a \ b \ \text{false}) \ \text{false} \ t \\ &\rightsquigarrow^* \text{false } t \ \text{false} = (\lambda x. \lambda y. y) \ t \ \text{false} \rightsquigarrow^* \text{false} \end{aligned}$$

All the other logical operators could be defined if we had a conditional construct *if-then-else*. In fact this can be defined very easily:

$$\text{if} := \lambda b. \lambda u. \lambda v. b \ u \ v.$$

Let's verify that it has the correct computational behaviour:

$$\begin{aligned} \text{if true } t_1 \ t_2 &= (\lambda b. \lambda u. \lambda v. b \ u \ v) \ \text{true} \ t_1 \ t_2 \\ &\rightsquigarrow^* \text{true } t_1 \ t_2 = (\lambda x. \lambda y. x) \ t_1 \ t_2 \rightsquigarrow^* t_1 \\ \text{if false } t_1 \ t_2 &= (\lambda b. \lambda u. \lambda v. b \ u \ v) \ \text{false} \ t_1 \ t_2 \\ &\rightsquigarrow^* \text{false } t_1 \ t_2 = (\lambda x. \lambda y. y) \ t_1 \ t_2 \rightsquigarrow^* t_2. \end{aligned}$$

Then we can define all logical connectives as conditionals, for example:

$$\text{and} := \lambda a. \lambda b. \text{if } a \ b \ \text{false}, \quad \text{or} := \lambda a. \lambda b. \text{if } a \ \text{true } b, \quad \text{not} := \lambda a. \text{if } a \ \text{false } \text{true}.$$

As a curiosity (just a feature of the encoding, don't read anything deeply philosophical in it) notice that *false is the true identity!*

$$\text{true id} \rightsquigarrow \text{false}$$

This is also a good example in the managing of variables:  $\text{true} = \lambda x. \lambda y. x$  and  $\text{id} = \lambda x. x$ . If we follow the variable convention, we shouldn't use  $x$  twice, so let's rename it in the identity:  $\text{id} = \lambda z. z$ . Then

$$\text{true id} = (\lambda x. \lambda y. x) (\lambda z. z) \rightsquigarrow \lambda y. \lambda z. z = \lambda x. \lambda y. y = \text{false}$$

where at the end we're free to rename the bound variables according to  $\alpha$ -conversion.

**Tuples** Pairs of  $\lambda$ -terms can be encoded by a single term: If  $t_1$  and  $t_2$  are terms, we define the encoding of the pair as

$$\langle t_1, t_2 \rangle := \lambda x. x t_1 t_2.$$

First and second projections are obtained by applying a pair to the familiar projections (or truth values) that we have already seen:

$$\begin{aligned} \text{fst } p &= p (\lambda x. \lambda y. x) \\ \text{snd } p &= p (\lambda x. \lambda y. y) \end{aligned}$$

We can verify that they have the correct reduction behaviour:

$$\begin{aligned} \text{fst } \langle t_1, t_2 \rangle &= \langle t_1, t_2 \rangle (\lambda x. \lambda y. x) = (\lambda x. x t_1 t_2) (\lambda x. \lambda y. x) \rightsquigarrow (\lambda x. \lambda y. x) t_1 t_2 \rightsquigarrow^* t_1, \\ \text{snd } \langle t_1, t_2 \rangle &= \langle t_1, t_2 \rangle (\lambda x. \lambda y. y) = (\lambda x. x t_1 t_2) (\lambda x. \lambda y. y) \rightsquigarrow (\lambda x. \lambda y. y) t_1 t_2 \rightsquigarrow^* t_2. \end{aligned}$$

Triples and longer tuples may be encoded as repeated pairs, for example  $\langle t_1, t_2, t_3 \rangle := \langle t_1, \langle t_2, t_3 \rangle \rangle$ , or directly using the same idea as for pairs:  $\langle t_1, t_2, t_3 \rangle := \lambda x. x t_1 t_2 t_3$ .

## 12.4 Confluence

A  $\lambda$ -term may contain several redexes. We have the choice of which one to reduce first. When we make one step of  $\beta$ -reduction, some of the redexes that were there in the beginning may disappear, some may be duplicated into many copies, some new ones may be created. Although we say that  $\beta$ -reduction is a “simplification” of the term, in the sense that we eliminate a pair of consecutive abstraction and application by immediately performing the associated substitution, the resulting reduced term is not always simpler. It may actually be much longer and more complicated.

Therefore it is not obvious that, if we choose different redexes to simplify, we will eventually get the same result. It is also not clear whether the reduction of a term will eventually terminate.

The first property is anyway true. But there are indeed terms whose reduction does not terminate.

**Theorem 12.1 (Confluence)** *Given any  $\lambda$ -term  $t$ , if  $t_1$  and  $t_2$  are two reducts of it, that is  $t \rightsquigarrow^* t_1$  and  $t \rightsquigarrow^* t_2$ ; then there exists a common reduct  $t_3$  such that  $t_1 \rightsquigarrow^* t_3$  and  $t_2 \rightsquigarrow^* t_3$ .*

**Proof.** TO BE COMPLETED □

**Definition 12.2** A normal form is a  $\lambda$ -term that doesn't contain any redexes. A term weakly normalizes if there is a sequence of reduction steps that ends in a normal form. A term strongly normalizes if any sequence of reduction steps eventually ends in a normal form.

There exist terms that do not normalize. The most famous one is a very short expression that reduces to itself:

$$\begin{aligned}\omega &:= (\lambda x.xx)(\lambda x.xx) \\ &\rightsquigarrow (xx)[x := \lambda x.xx] = (\lambda x.xx)(\lambda x.xx) = \omega \\ &\rightsquigarrow \omega \rightsquigarrow \dots\end{aligned}$$

There are also terms that grow without bound when we reduce them, for example:

$$\begin{aligned}(\lambda x.xx x)(\lambda x.xx x) &\rightsquigarrow (\lambda x.xx x)(\lambda x.xx x)(\lambda x.xx x) \\ &\rightsquigarrow (\lambda x.xx x)(\lambda x.xx x)(\lambda x.xx x)(\lambda x.xx x) \\ &\rightsquigarrow \dots\end{aligned}$$

Here is an example of a term that weakly normalizes but doesn't strongly normalize:

$$(\lambda x.\lambda y.x)(\lambda z.z)\omega.$$

This term applies the first projection function to two arguments. If we immediately reduce the application of the projection, the argument  $\omega$  disappears and we are left with the identity function, which is a normal form:

$$(\lambda x.\lambda y.x)(\lambda z.z)\omega \rightsquigarrow \lambda z.z.$$

However, if we try to reduce the redex inside the second argument, we don't make any progress and we could continue reducing it forever.

## 12.5 Recursion

We have seen how to define some basic arithmetic functions as  $\lambda$ -terms: addition `plus`, multiplication `mult`, exponentiation `exp`. With a little effort you may be able to define the predecessor function and subtraction. What about (whole) division, maybe using the Euclid algorithm?

If we want to use the  $\lambda$ -calculus as a complete programming language, there should be a way to define any computable function. This is indeed possible. In fact there is a single  $\lambda$ -term, called *the Y combinator*, that allows us to use unrestricted recursion:

$$Y := \lambda f.(\lambda x.f(x x))(\lambda x.f(x x)).$$

This definition is inspired by the non-normalizing term  $\omega$ . In fact we have that  $Y \text{ id} \rightsquigarrow^* \omega$ .

The most striking property of  $Y$  is that it computes a *fixed point* for every term  $F$ :

$$Y F \rightsquigarrow^* F(Y F).$$

This is not exactly true: to be precise,  $Y F \rightsquigarrow (\lambda x. F(x x)) (\lambda x. F(x x)) =: \text{fix}_F$  and  $\text{fix}_F \rightsquigarrow^* F(\text{fix}_F)$ . If we keep reducing, we get an infinite sequence of applications of  $F$ :

$$\text{fix}_F \rightsquigarrow^* F \text{fix}_F \rightsquigarrow^* F(F \text{fix}_F) \rightsquigarrow^* F(F(F \text{fix}_F)) \rightsquigarrow^* \dots$$

To define a recursive function, we just need to encode a single step of it as a term  $F$  and then use  $Y$  to iterate that step as many times as it is necessary to get a result.

Here is, for example, the definition of the factorial:

$$\begin{aligned} \text{fact}_{\text{step}} &:= \lambda f. \lambda n. \text{if } (\text{isZero } n) \bar{1} (\text{mult } n (f(\text{pred } n))) \\ \text{fact} &:= Y \text{fact}_{\text{step}}. \end{aligned}$$

## 12.6 The universality of $\lambda$ -calculus

The  $\lambda$ -calculus was invented to answer the question: “What does it mean that a problem is effectively solvable?” Up to that point the notion of a question being answerable by an precise method was left to intuition. Throughout the history of mathematics many difficult problems were posed. When someone claimed to have a solution, expert mathematicians would analyze it and come to an informed option about whether it was correct and precise.

But at the beginning of the 20th century, David Hilbert had formulated the challenge of defining exactly what we mean by a precise method: can we give a mathematically rigorous definition of what an effective procedure is? Both Alonzo Church and Alan Turing worked towards a realization of Hilbert’s challenge and they both came up with a fully satisfying solution.

But their two theories looked completely different. Which one would be the “true” answer to Hilbert’s question? It turned out that, beyond their formal difference, Turing machines and  $\lambda$ -terms are equally expressive. The same set of computable functions can be implemented in both. Turing himself proved this when he discovered Church’s work and added a sketch of the proof as an Appendix to his article about computable numbers.

Around the same time, other models of computation were proposed. Kurt Gödel formulated the notion of  $\mu$ -recursive function: any function on the natural numbers defined by certain forms of recursive definition. Stephen Kleene, a student of Church, proved that these functions are exactly those definable in the  $\lambda$ -calculus and went on to develop a rich theory of computation based on them.

Another of Church’s students, J. Barkley Rosser, was the first to clearly formulate the notion that these three definitions were equivalent realizations of the informal notion of an effective computation method: “All three definitions are equivalent, so it does not matter which one is used.”

Later, when the first digital computers were constructed, the Hungarian-American mathematician John von Neumann proposed a model that is more realistic and describes the actual architecture of real computers: the *register machine*. This notion was also equivalent to the previous ones. Many other different characterizations have been proposed since and they all turned out to be equivalent.

The statement that every effective definition of computable processes is equivalent to the  $\lambda$ -calculus (or to Turing machines or to  $\mu$ -recursive functions



and so on) is known as *the Church-Turing Thesis* (Stephen Kleene was the one who actually first formulated and named it).

The  $\lambda$ -calculus is not just a mathematical abstraction. All functional programming languages are based on (a typed version of) it. Lisp, ML, Haskell, OCaml are the most well-known. Functional programming has been so successful that traditional imperative languages are starting to introduce function abstraction as a basic feature in their definition. For example JavaScript and Python contain first-class functions as part of their definition.

## 12.7 Exercises

### Exercise 12.1

Consider the following  $\lambda$ -terms: `nand-pair` is a function on pairs of Booleans, `nand-fun` a function from Church Numerals to pairs of Booleans:

$$\begin{aligned}\text{nand-pair} &= \lambda p. \langle \text{not } (\text{and } (p \text{ true}) (p \text{ false})), p \text{ true} \rangle \\ \text{nand-fun} &= \lambda n. n \text{ nand-pair } \langle \text{false}, \text{false} \rangle\end{aligned}$$

1. What values do the following terms reduce to?

$$\begin{array}{ll}\text{nand-pair } \langle \text{true}, \text{true} \rangle \rightsquigarrow^* ? & \text{nand-pair } \langle \text{false}, \text{true} \rangle \rightsquigarrow^* ? \\ \text{nand-pair } \langle \text{true}, \text{false} \rangle \rightsquigarrow^* ? & \text{nand-pair } \langle \text{false}, \text{false} \rangle \rightsquigarrow^* ?\end{array}$$

2. Show the steps of reduction in the computation of  $(\text{nand-fun } \bar{4})$ . [You can use the previous reductions and those from the lecture notes as single steps.]
3. Give an informal definition of what `nand-fun` does: for which numbers  $n$  does  $(\text{nand-fun } \bar{n}) \rightsquigarrow^* \langle \text{true}, \text{true} \rangle$ ?

### Exercise 12.2

Write a  $\lambda$ -term that implements the following function:

$$\begin{aligned}\text{thrFib} &: \mathbb{N} \rightarrow \mathbb{N} \\ \text{thrFib } 0 &= 0 \\ \text{thrFib } 1 &= 0 \\ \text{thrFib } 2 &= 1 \\ \text{thrFib } (n + 3) &= \text{thrFib } n + \text{thrFib } (n + 1) + \text{thrFib } (n + 2)\end{aligned}$$

[Use the same idea that was used in the lecture to define the Fibonacci numbers: first write an auxiliary function that returns a triple of numbers and then extract the first.]

## 13 Algorithmic Complexity

The undecidability of the Halting Problem shows that there are tasks that are impossible to accomplish using any computing machine. This is an intrinsic limitation of computation. But even for some problems for which algorithmic solutions are known, it may be impossible to compute them effectively. In some cases, the search for a solution is so complex that it would take a time longer than the age of the universe to get the answer. We can always hope for dramatic improvements in the power and speed of computers, but some tasks seem to have an essentially intractable complexity.

In talking about the time complexity of algorithms, we want to abstract away from implementation details. The exact number of steps needed to compute a result may depend on the model of computation: the same algorithm, realized as a Turing Machine or as a  $\lambda$ -term, will result in different number of head operations of the machine and reduction steps of the  $\lambda$ -term. The exact time of the computation will depend on the specific machine on which it runs: even the exact same program in the exact same language will take different times on different computers.

But beyond these differing details, the *complexity classes* of algorithms are quite constant across different models of computation and different computing devices. A complexity class characterizes algorithms that have a relation between size of the input and time of computation that can be expressed by certain functions. Although they are traditionally formulated in terms of Turing Machines, we can show that realizations in different paradigms have similar behaviour.

The most important complexity classes are  $\mathcal{P}$  and  $\mathcal{NP}$ .  $\mathcal{P}$  stands for *Polynomial Complexity*: A problem is in  $\mathcal{P}$  if it can be solved by a Turing Machine whose running time is at most a polynomial in the size of the input.  $\mathcal{NP}$  stands for *Non-deterministic Polynomial Complexity*: A problem is in  $\mathcal{NP}$  if it can be solved by a Non-deterministic Turing Machine (NTM) whose running time is at most a polynomial in the size of the input. Remember that an NTM may have overlapping transitions: in certain configurations, the machine may have several allowed steps. In every run, it will randomly choose which transition to apply. A word is accepted if there exists one possible computation (among all allowed by the non-deterministic choices) that reaches an accepting state.

There are two alternative and equivalent ways of seeing the class  $\mathcal{NP}$ . A problem is in  $\mathcal{NP}$  if it is solvable in polynomial time by a *parallel computer* which is allowed to spawn several parallel computations. Each computation is independent and must terminate in polynomial time; it is sufficient that the solution is found by one of the computations. Although the running time is polynomial, there may be an exponential blow-up of the number of parallel computations that need to run simultaneously.

The second alternative explanation of the class  $\mathcal{NP}$  is that a problem is in it if we can verify solutions in polynomial time. There is an algorithm that, when given as input an instance of the problem and a potential solution, terminates in polynomial time and gives a positive answer if and only if the solution is correct.

Intuitively, it seems obvious that a problem that belongs to  $\mathcal{NP}$  doesn't necessarily belong to  $\mathcal{P}$ : being able to find a solution in polynomial time by using non-deterministic or parallel computations, or being able to verify a solution in

polynomial time, doesn't mean that we can generate a solution deterministically in polynomial time. There may be an exponential number of non-deterministic or parallel computations and an exponential number of potential solutions to verify.

However, until now nobody managed to prove that  $\mathcal{P} \neq \mathcal{NP}$ . It is the most famous open problem in theoretical computer science. There are strong reasons to believe that the two classes are distinct. There is a subclass of  $\mathcal{NP}$ , called  $\mathcal{NP}$ -complete, that contains problems that are in a sense universal with respect to this issue. A problem is in  $\mathcal{NP}$ -complete if it is in  $\mathcal{NP}$  and every other problem in  $\mathcal{NP}$  can be reduced to it in polynomial time. This means that if somebody finds a polynomial-time algorithm to solve one of these problems, then automatically all  $\mathcal{NP}$  problems will be solvable in polynomial time. To date, thousands of different problems from disparate branches of computer science and mathematics have been proved to be  $\mathcal{NP}$ -complete. It seems highly unlikely that we could solve them all at once with the same algorithm.

### 13.1 The Satisfiability Problem

The first problem to be proved  $\mathcal{NP}$ -complete was Boolean Satisfiability (in short SAT). An instance of the problem is a propositional formula, that is, an expression that uses variables, the propositional connectives for conjunction ( $\wedge$ ), disjunction ( $\vee$ ) and negation ( $\neg$ ), and parentheses. An example is the formula  $(x_1 \vee \neg x_2) \wedge \neg x_1$ .

Solving an instance of SAT means determining if there is an assignment of truth values to the variables that makes the formula true. In the example above the assignment  $[x_1 \mapsto \text{false}, x_2 \mapsto \text{false}]$  is such a solution. Now consider the following more complicated formula in three variables:

$$f = (\neg x_1 \vee x_2) \wedge (x_3 \vee \neg x_2) \wedge \neg(x_3 \wedge \neg x_1) \wedge (x_1 \wedge \neg x_2 \vee x_2 \wedge \neg x_3).$$

There is no assignment of truth values to  $x_1, x_2, x_3$  that makes  $f$  true. To check this, we must try all potential solutions:

$$\begin{aligned} [x_1 \mapsto \text{true}, x_2 \mapsto \text{true}, x_3 \mapsto \text{true}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{true}, x_2 \mapsto \text{true}, x_3 \mapsto \text{false}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{true}, x_2 \mapsto \text{false}, x_3 \mapsto \text{true}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{true}, x_2 \mapsto \text{false}, x_3 \mapsto \text{false}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{false}, x_2 \mapsto \text{true}, x_3 \mapsto \text{true}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{false}, x_2 \mapsto \text{true}, x_3 \mapsto \text{false}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{false}, x_2 \mapsto \text{false}, x_3 \mapsto \text{true}] &\implies f \mapsto \text{false} \\ [x_1 \mapsto \text{false}, x_2 \mapsto \text{false}, x_3 \mapsto \text{false}] &\implies f \mapsto \text{false} \end{aligned}$$

In general, to look for a solution or to verify that there is no solution, we need to check all the possible assignments of truth values to the variables. If there are  $n$  variables in the formulas, then we must check  $2^n$  assignments. Therefore checking them all takes a time that is exponential in the number of variables, so exponential in the size of the input. On the other hand, it is straightforward to verify whether an assignment gives a solution or not.

### 13.2 Time Complexity

Let's define time complexity exactly using Turing Machines.

**Definition 13.1** A Turing Machine  $M$  is said to have time complexity  $f$ , a function  $\mathbb{N} \rightarrow \mathbb{N}$ , if, whenever we run  $M$  on an input of size  $n$ , it will take at most  $f(n)$  steps to terminate.

The number of steps is measured by counting the movement of the reading/writing head of  $M$ , as given by the next-state relation  $\vdash_M$ .

This is a *worst case* measure of the running time of the machine: It is possible that  $M$  will terminate in less than  $n$  steps on some (or all) inputs. So  $f$  provides just an upper bound. We are interested in getting the answer within a certain time constraint.

As we know, any “problem” can be expressed precisely as a language recognition task. Each instance of the problem is specified by giving a description in the form of a word/list of symbols. So defining complexity classes for problems is the same as defining them for languages.

**Definition 13.2** A language  $L$  is in the class  $\mathcal{P}$  if there is a Turing Machine  $M$  that decides  $L$  and has a time complexity expressible by a polynomial function.

In the previous two definitions, we were talking about *deterministic* Turing machines. In any given instantaneous description of  $M$ , the transition function determines uniquely the step to be performed.

We are also going to consider *non-deterministic* Turing machines. These may have several possible transitions in the same instantaneous description. The transition function has the type:

$$\delta \in Q \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma \times \{L, R\})$$

So, if  $M$  is in state  $q$  and the head is reading a symbol  $x$ , the transition  $\delta(q, x)$  is a set of triples, each one specifying a different action. For example, if  $\delta(q, x) = \{(q_1, y, L), (q_2, z, R)\}$ , then the machine may either write a  $y$ , move left and go into state  $q_1$  or write a  $z$ , move right and go into state  $q_2$ . When we run the machine, in these situation it will randomly choose which of the several possible steps to perform. We say that such a machine accepts a word if there is one run (among the many possible) that leads to an accepting state.

**Definition 13.3** A language  $L$  is in the class  $\mathcal{NP}$  if there is a non-deterministic Turing Machine  $M$  that decides  $L$  and has a time complexity expressible by a polynomial function. In this case the polynomial bound must apply to all possible runs of  $M$  on a given input word.

Equivalently, we can express membership in the class  $\mathcal{NP}$  using deterministic Turing machines that take as input both a word and a certificate that proves that the word belongs to  $L$ . A certificate can be any string of symbols that the machine may interpret as evidence.

**Theorem 13.4** A language  $L$  is in the class  $\mathcal{NP}$  if there is a deterministic Turing Machine  $M$  and a polynomial function  $f$  with the following properties:

- When we run  $M$  on a tape containing two inputs  $w$  and  $v$  (separated by a blank), it always terminates in a number of steps smaller than  $f(|w|)$ ;
- If  $w \in L$ , then there exists some  $v$  such that  $M$  gives a positive answer when run on  $w$  and  $v$ ;

- If  $w \notin L$ , then for every  $v$ ,  $M$  always gives a negative answer when run on  $w$  and  $v$ .

A string  $v$  such that  $M$  gives a positive answer when run on  $w$  and  $v$  is a *certificate* or proof that  $w$  is in  $L$ . If you think of  $L$  as representing some problem and  $w$  as an instance of it, then  $v$  can be thought of as a solution for  $w$ .

For example, given an instance of SAT, a certificate is an assignment of truth values to the variables that makes the formula true. We can define a Turing Machine with polynomial time complexity that, when run on a formula and an assignment of values to the variables, determines whether the formula computes to true under that assignment.

**Theorem 13.5** SAT is in  $\mathcal{NP}$ .

### 13.3 $\mathcal{NP}$ -completeness

Among the  $\mathcal{NP}$  problems there are some that are *universal* in the sense that all  $\mathcal{NP}$  problems can be reduced to them in polynomial time.

**Definition 13.6** A language  $L_0$  is in  $\mathcal{NP}$ -complete if it is  $\mathcal{NP}$  and, for every language  $L_1$  in  $\mathcal{NP}$  there exists a Turing Machine  $M$  that runs in polynomial time, with the following property:

When we run  $M$  on a word  $w_1$ , it will terminate with the tape containing a word  $w_0$  such that  $w_1 \in L_1$  if and only if  $w_0 \in L_0$ .

We say that  $L_1$  is reduced to  $L_0$  in polynomial time.

Seeing  $L_0$  and  $L_1$  as encodings of problems  $P_0$  and  $P_1$ , the definition says that we can turn every problem of  $P_1$  into a problem of  $P_0$  with the same solution.

The great turning point in the study of the time complexity of algorithms came with the discovery that there exist some  $\mathcal{NP}$ -complete problems.

**Theorem 13.7 (Stephen Cook, 1971)** SAT is  $\mathcal{NP}$ -complete.

Because Cook's Theorem, many other problems have been shown to be  $\mathcal{NP}$ -complete. They now run in the thousands. Here's a brief description of some famous ones.

**THE TRAVELLING SALESMAN** A salesman must visit  $n$  different towns. He has a table giving the distances between any pair of towns. He has a budget that he can use to buy petrol, which will allow him to travel a maximum length. Is there a route around all the towns with length shorter or equal to that allowed by the budget?

**SUBSET SUM** Given a set of (positive and negative) integers, is there a non-empty subset of it that sums up to zero? This is a special case of the *knapsack* problem that consists in maximizing the value of a set of objects that can fit into a knapsack with a limited capacity.

**GRAPH COLOURING** Given a graph and a fixed number of colours, is it possible to colour all the nodes of the graph so that two nodes of the same colour are not linked by an edge?

To this day, it is still unknown whether there are  $\mathcal{NP}$  problems that are not in  $\mathcal{P}$ . All attempts to find a polynomial algorithm that solves any  $\mathcal{NP}$ -complete problem have failed. Also all attempts to prove that such an algorithm doesn't exist have failed.

If we were to discover such an algorithm, it would automatically give us a way to solve all  $\mathcal{NP}$ -complete problems in polynomial time: we can just reduce them all to the one we can solve. For this reason, most mathematicians and computer scientists think that it is impossible, but the question is still open and may be *the greatest mystery in computer science*:

$$\mathcal{P} = \mathcal{NP} \quad ?$$

## 13.4 Exercises

### Exercise 13.1

This and the next exercise concerns programing a SAT solver in Haskell. The first exercise is to write a program that *checks* whether a proposed solution to an instance of SAT is correct. Use the following type definitions:

```
data SAT = Var Int | Not SAT | And SAT SAT | Or SAT SAT

type Assignment = [Bool]
```

SAT represents Boolean formulas with variables `Var 0`, `Var 1`, `Var 2`, `Var 3` and so on. An assignment is a list of Booleans, giving the values to some of the variables. For example `[True,False,False,True]` assigns the values: `Var 0 = True`, `Var 1 = False`, `Var 2 = False`, `Var 3 = True`. The other variables are left without value.

Write an evaluation function:

```
evaluate :: SAT -> Assignment -> Bool
```

It uses an assignment of values to variables to evaluate a formula. (If the formula contains variables with indices larger or equal to the length of the assignment, you can leave it undefined.)

### Exercise 13.2

This exercise concerns writing a program that *decides* the solvability of an instance of SAT. We will break the task down into several steps.

1. Write a function that gives the highest index of a variable occurring in a formula:

```
varNum :: SAT -> Int
```

2. Write a function that generates all the assignments for all variables up to a given index:

```
allAssign :: Int -> [Assignment]
```

For example `allAssign 2` should give all the possible assignments for variables `Var 0`, `Var 1` and `Var 2`:

```
allAssign 2 = [[True,True,True],   [True,True,False],
               [True,False,True],  [True,False,False],
               [False,True,True],   [False,True,False],
               [False,False,True],  [False,False,False]]
```

3. Write a function that, for a given formula, verifies if there is one assignment on which the formula evaluates to `True`:

```
satisfiable :: SAT -> Bool
```

4. Define a function that actually returns the solution, if it exists:

```
solution :: SAT -> Maybe Assignment
```

### Exercise 13.3

Discuss informally the complexity properties of the programs that you wrote.

1. Give an informal description of the steps of computation of `evaluate` and explain why it runs in polynomial time on the length of the input. Use this to argue that SAT is in the class  $\mathcal{NP}$ .
2. Give an informal description of the steps of computation of `satisfiable` and explain why it runs in exponential time on the length of the input. *(If you managed to write a polynomial-time program, congratulations: you solved  $\mathcal{P} = \mathcal{NP}$ !)*
3. Explain what it means that SAT is  $\mathcal{NP}$ -complete and what relevance this has for the  $\mathcal{P} = \mathcal{NP}$  question.

## References

- [ASU86] Alfred V. Aho, Ravi Sethi, and Ullman Jeffrey D. Ullman. *Compilers - Principles, Techniques, and Tools*. Addison-Wesley, 1986. 87
- [Cox07] Russ Cox. Regular expression matching can be simple and fast (but is slow in java, perl, php, python, ruby, ...). <http://swtch.com/rsc/regexp/regexp1.html>, January 2007. 26
- [GJS<sup>+</sup>15] James Gosling, Bill Joy, Guy Steele, Gilad Bracha, and Alex Buckley. *The Java Language Specification*. Oracle, Inc., Java SE 8 edition edition, 2015. 59
- [HMU01] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison Wesley, 2nd edition edition, 2001. 14, 41, 76, 100
- [Hud99] Scott E. Hudson. Cup parser generator for java. <http://www.cs.princeton.edu/~appel/modern/java/CUP/><http://www.cs.princeton.edu/ap-pel/modern/java/CUP/>, 1999. 96
- [Mar01] Simon Marlow. Happy: The parser generator for haskell. <http://www.haskell.org/happy/><http://www.haskell.org/happy/>, 2001. 96
- [Nil16] Henrik Nilsson. Compilers (G53CMP) — lecture notes. <http://www.cs.nott.ac.uk/~nhn/G53CMP><http://www.cs.nott.ac.uk/~nhn/G53CMP>, 2016. 59, 79
- [Par05] Terence Parr. ANTLR: Another tool for language recognition. <http://wwwantlr.org/><http://wwwantlr.org/>, 2005. 96