

Languages and Compilers

Johan Jeuring, Doaitse Swierstra

January 29, 2017

Copyright © 2001–2009 by Johan Jeuring and Doaitse Swierstra.

Contents

Preface for the 2009 version	v
Voorwoord	vii
1. Goals	1
1.1. History	2
1.2. Grammar analysis of context-free grammars	3
1.3. Compositionality	4
1.4. Abstraction mechanisms	4
2. Context-Free Grammars	7
2.1. Languages	8
2.2. Grammars	11
2.2.1. Notational conventions	14
2.3. The language of a grammar	15
2.3.1. Examples of basic languages	17
2.4. Parse trees	19
2.5. Grammar transformations	22
2.5.1. Removing duplicate productions	23
2.5.2. Substituting right hand sides for nonterminals	23
2.5.3. Removing unreachable productions	23
2.5.4. Left factoring	24
2.5.5. Removing left recursion	24
2.5.6. Associative separator	25
2.5.7. Introduction of priorities	26
2.5.8. Discussion	27
2.6. Concrete and abstract syntax	28
2.7. Constructions on grammars	31
2.7.1. SL: an example	34
2.8. Parsing	36
2.9. Exercises	37
3. Parser combinators	41
3.1. The type of parsers	43
3.2. Elementary parsers	45
3.3. Parser combinators	48

3.3.1. Matching parentheses: an example	53
3.4. More parser combinators	55
3.4.1. Parser combinators for EBNF	55
3.4.2. Separators	58
3.5. Arithmetical expressions	60
3.6. Generalised expressions	63
3.7. Exercises	64
4. Grammar and Parser design	67
4.1. Step 1: Example sentences for the language	68
4.2. Step 2: A grammar for the language	68
4.3. Step 3: Testing the grammar	69
4.4. Step 4: Analysing the grammar	69
4.5. Step 5: Transforming the grammar	69
4.6. Step 6: Deciding on the types	70
4.7. Step 7: Constructing the basic parser	71
4.7.1. Basic parsers from strings	72
4.7.2. A basic parser from tokens	72
4.8. Step 8: Adding semantic functions	74
4.9. Step 9: Did you get what you expected	75
4.10. Exercises	76
5. Compositionality	77
5.1. Lists	78
5.1.1. Built-in lists	78
5.1.2. User-defined lists	80
5.1.3. Streams	81
5.2. Trees	82
5.2.1. Binary trees	82
5.2.2. Trees for matching parentheses	83
5.2.3. Expression trees	86
5.2.4. General trees	88
5.2.5. Efficiency	88
5.3. Algebraic semantics	89
5.4. Expressions	90
5.4.1. Evaluating expressions	90
5.4.2. Adding variables	91
5.4.3. Adding definitions	93
5.4.4. Compiling to a stack machine	95
5.5. Block structured languages	97
5.5.1. Blocks	97
5.5.2. Generating code	98
5.6. Exercises	102

6. Computing with parsers	105
6.1. Insert a semantic function in the parser	105
6.2. Apply a fold to the abstract syntax	105
6.3. Deforestation	106
6.4. Using a class instead of abstract syntax	107
6.5. Passing an algebra to the parser	109
7. Programming with higher-order folds	111
7.1. The rep_min problem	112
7.1.1. A straightforward solution	113
7.1.2. Lambda lifting	113
7.1.3. Tupling computations	114
7.1.4. Merging tupled functions	115
7.2. A small compiler	117
7.2.1. The language	117
7.2.2. A stack machine	117
7.2.3. Compiling to the stackmachine	118
7.3. Attribute grammars	122
8. Regular Languages	123
8.1. Finite-state automata	123
8.1.1. Deterministic finite-state automata	124
8.1.2. Nondeterministic finite-state automata	126
8.1.3. Implementation	129
8.1.4. Constructing a DFA from an NFA	132
8.1.5. Partial evaluation of NFA's	135
8.2. Regular grammars	136
8.2.1. Equivalence of Regular grammars and Finite automata	139
8.3. Regular expressions	142
8.4. Proofs	145
8.4.1. Proof of Theorem 8.7	146
8.4.2. Proof of Theorem 8.16	147
8.4.3. Proof of Theorem 8.17	148
8.5. Exercises	150
9. Pumping Lemmas: the expressive power of languages	153
9.1. The Chomsky hierarchy	154
9.1.1. Type-0 grammars	154
9.1.2. Type-1 grammars	154
9.1.3. Type-2 grammars	154
9.1.4. Type-3 grammars	155
9.2. The pumping lemma for regular languages	155
9.3. The pumping lemma for context-free languages	157
9.4. Proofs of pumping lemmas	161

9.4.1. Proof of the Regular Pumping Lemma, Theorem 9.1	161
9.4.2. Proof of the Context-free Pumping Lemma, Theorem 9.2	162
9.5. Exercises	164
10.LL Parsing	167
10.1. LL Parsing: Background	167
10.1.1. A stack machine for parsing	168
10.1.2. Some example derivations	169
10.1.3. <i>LL</i> (1) grammars	173
10.2. LL Parsing: Implementation	176
10.2.1. Context-free grammars in Haskell	176
10.2.2. Parse trees in Haskell	178
10.2.3. <i>LL</i> (1) parsing	179
10.2.4. Implementation of <i>isLL</i> (1)	181
10.2.5. Implementation of lookahead	181
10.2.6. Implementation of empty	186
10.2.7. Implementation of first and last	187
10.2.8. Implementation of follow	190
11.LL versus LR parsing	195
11.1. <i>LL</i> (1) parser example	196
11.1.1. An <i>LL</i> (1) checker	196
11.1.2. An <i>LL</i> (1) grammar	197
11.1.3. Using the <i>LL</i> (1) parser	198
11.2. LR parsing	199
11.2.1. A stack machine for <i>SLR</i> parsing	200
11.3. LR parse example	201
11.3.1. An LR checker	201
11.3.2. LR action selection	202
11.3.3. LR optimizations and generalizations	206
A. The Stack module	211
B. Answers to exercises	213

Preface for the 2009 version

This is a work in progress. These lecture notes are in large parts identical with the old lecture notes for “Grammars and Parsing” by Johan Jeuring and Doaitse Swierstra.

The lecture notes have not only been used at Utrecht University, but also at the Open University. Some modifications made by Manuela Witsiers have been reintegrated.

I have also started to make some modifications to style and content – mainly trying to adapt the Haskell code contained in the lecture notes to currently established coding guidelines. I also rearranged some of the material because I think they connect better in the new order.

However, this work is not quite finished, and this means that some of the later chapters look a bit different from the earlier chapters. I apologize in advance for any inconsistencies or mistakes I may have introduced.

Please feel free to point out any mistakes you find in these lecture notes to me – any other feedback is also welcome.

I hope to be able to update the online version of the lecture notes regularly.

Andres Löh

October 2009

Preface for the 2009 version

Voorwoord

Het hiervolgende dictaat is gebaseerd op teksten uit vorige jaren, die onder andere geschreven zijn in het kader van het project *Kwaliteit en Studeerbaarheid*.

Het dictaat is de afgelopen jaren verbeterd, maar we houden ons van harte aanbevolen voor suggesties voor verdere verbetering, met name daar waar het het aangeven van verbanden met andere vakken betreft.

Veel mensen hebben een bijgedrage geleverd aan de totstandkoming van dit dictaat door een gedeelte te schrijven, of (een gedeelte van) het dictaat te becommentariëren. Speciale vermelding verdienen Jeroen Fokker, Rik van Geldrop, en Luc Duponcheel, die mee hebben geholpen door het schrijven van (een) hoofdstuk(ken) van het dictaat. Commentaar is onder andere geleverd door: Arthur Baars, Arnoud Berendsen, Gijsbert Bol, Breght Boschker, Martin Bravenboer, Pieter Eendebak, Alexander Elyasov, Matthias Felleisen, Rijk-Jan van Haften, Graham Hutton, Daan Leijen, Andres Löh, Erik Meijer, en Vincent Oostindië.

Tenslotte willen we van de gelegenheid gebruik maken enige *studeeraanwijzingen* te geven:

- Het is onze eigen ervaring dat het uitleggen van de stof aan iemand anders vaak pas duidelijk maakt welke onderdelen je zelf nog niet goed beheerst. Als je dus van mening bent dat je een hoofdstuk goed begrijpt, probeer dan eens *in eigen woorden uiteen te zetten*.
- *Oefening baart kunst*. Naarmate er meer aandacht wordt besteed aan de presentatie van de stof, en naarmate er meer voorbeelden gegeven worden, is het verleidelijker om, na lezing van een hoofdstuk, de conclusie te trekken dat je een en ander daadwerkelijk beheerst. “Begrijpen is echter niet hetzelfde als “kennen”, “kennen” is iets anders dan “beheersen” en “beheersen” is weer iets anders dan “er iets mee kunnen”. Maak dus de opgaven die in het dictaat opgenomen zijn zelf, en doe dat niet door te kijken of je de oplossingen die anderen gevonden hebben, begrijpt. Probeer voor jezelf bij te houden welk stadium je bereikt hebt met betrekking tot alle genoemde leerdoelen. In het ideale geval zou je in staat moeten zijn een mooi tentamen in elkaar te zetten voor je mede-studenten!
- *Zorg dat je up-to-date bent*. In tegenstelling tot sommige andere vakken is het bij dit vak gemakkelijk de vaste grond onder je voeten kwijt te raken. Het is niet “elke week nieuwe kansen”. We hebben geprobeerd door de indeling van

Voorwoord

de stof hier wel iets aan te doen, maar de totale opbouw laat hier niet heel veel vrijheid toe. Als je een week gemist hebt is het vrijwel onmogelijk de nieuwe stof van de week daarop te begrijpen. De tijd die je dan op college en werkcollege doorbrengt is dan weinig effectief, met als gevolg dat je vaak voor het tentamen heel veel tijd (die er dan niet is) kwijt bent om in je uppie alles te bestuderen.

- We maken gebruik van de taal Haskell om veel concepten en algoritmen te presenteren. Als je nog moeilijkheden hebt met de taal Haskell aarzel dan niet direct hier wat aan te doen, en zonodig hulp te vragen. Anders maak je jezelf het leven heel moeilijk. Goed gereedschap is het halve werk, en Haskell is hier ons gereedschap.

Veel sterkte, en hopelijk ook veel plezier,

Johan Jeuring en Doaitse Swierstra

1. Goals

Introduction

Courses on *Grammars*, *Parsing* and *Compilation of programming languages* have always been some of the core components of a computer science curriculum. The reason for this is that from the very beginning of these curricula it has been one of the few areas where the development of formal methods and the application of formal techniques in actual program construction come together. For a long time the construction of compilers has been one of the few areas where we had a methodology available, where we had tools for generating parts of compilers out of formal descriptions of the tasks to be performed, and where such program generators were indeed generating programs which would have been impossible to create by hand. For many practicing computer scientists the course on compiler construction still is one of the highlights of their education.

One of the things which were not so clear however is where exactly this joy originated from: the techniques taught definitely had a certain elegance, we could construct programs someone else could not – thus giving us the feeling we had “the right stuff” –, and when completing the practical exercises, which invariably consisted of constructing a compiler for some toy language, we had the usual satisfied feeling. This feeling was augmented by the fact that we would not have had the foggiest idea how to complete such a product a few months before, and now we knew “how to do it”.

This situation has remained so for years, and it is only in the last years that we have started to discover and make explicit the reasons why this area attracted so much interest. Many of the techniques which were taught on a “this is how you solve this kind of problems” basis, have been provided with a theoretical underpinning which explains why the techniques work. As a beneficial side-effect we also gradually learned to see where the discovered concept further played a rôle, thus linking the area with many other areas of computer science; and not only that, but also giving us a means to explain such links, stress their importance, show correspondences and transfer insights from one area of interest to the other.

Goals

The goals of these lecture notes can be split into primary goals, which are associated with the specific subject studied, and secondary – but not less important – goals which

1. Goals

have to do with developing skills which one would expect every educated computer scientist to have. The primary, somewhat more traditional, goals are to learn:

- to *describe* structures (i.e., “formulas”) using *grammars*;
- to *parse*, i.e., to recognise (build) such structures in (from) a sequence of symbols;
- to *analyse* grammars to see whether or not specific properties hold;
- to understand the concept of *compositionality*;
- to *apply these techniques* in the construction of all kinds of programs;
- to familiarise oneself with the concept of *computability*.

The secondary, more far reaching, goals are:

- to develop the capability *to abstract*;
- to understand the concepts of *abstract interpretation* and *partial evaluation*;
- to understand the concept of *domain specific languages*;
- to show how proper formalisations can be used as a starting point for the *construction of useful tools*;
- to improve the general *programming skills*;
- to show a wide variety of useful *programming techniques*;
- to show how to develop programs in a *calculational style*.

1.1. History

When at the end of the fifties the use of computers became more and more widespread, and their reliability had increased enough to justify applying them to a wide range of problems, it was no longer the actual hardware which posed most of the problems. Writing larger and larger programs by more and more people sparked the development of the first more or less machine-independent programming language FORTRAN (FORmula TRANslator), which was soon to be followed by ALGOL-60 and COBOL.

For the developers of the FORTRAN language, of which John Backus was the prime architect, the problem of how to describe the language was not a hot issue: much more important problems were to be solved, such as, what should be in the language and what not, how to construct a compiler for the language that would fit into the small memories which were available at that time (kilobytes instead of megabytes), and how to generate machine code that would not be ridiculed by programmers who had thus far written such code by hand. As a result the language was very much implicitly defined by what was accepted by the compiler and what not.

Soon after the development of FORTRAN an international working group started to work on the design of a machine independent high-level programming language, to become known under the name ALGOL-60. As a remarkable side-effect of this undertaking, and probably caused by the need to exchange proposals in writing, not only

a language standard was produced, but also a notation for describing programming languages was proposed by Naur and used to describe the language in the famous Algol-60 report. Ever since it was introduced, this notation, which soon became to be known as the Backus-Naur formalism (BNF), has been used as the primary tool for describing the basic structure of programming languages.

It was not for long that computer scientists, and especially people writing compilers, discovered that the formalism was not only useful to express what language should be accepted by their compilers, but could also be used as a guideline for structuring their compilers. Once this relationship between a piece of BNF and a compiler became well understood, programs emerged which take such a piece of language description as input, and produce a skeleton of the desired compiler. Such programs are now known under the name *parser generators*.

Besides these very mundane goals, i.e., the construction of compilers, the BNF-formalism also became soon a subject of study for the more theoretically oriented. It appeared that the BNF-formalism actually was a member of a hierarchy of *grammar classes* which had been formulated a number of years before by the linguist Noam Chomsky in an attempt to capture the concept of a “language”. Questions arose about the *expressibility* of BNF, i.e., which classes of languages can be expressed by means of BNF and which not, and consequently how to express restrictions and properties of languages for which the BNF-formalism is not powerful enough. In the lectures we will see many examples of this.

1.2. Grammar analysis of context-free grammars

Nowadays the use of the word Backus-Naur is gradually diminishing, and, inspired by the Chomsky hierarchy, we most often speak of *context-free grammars*. For the construction of everyday compilers for everyday languages it appears that this class is still a bit too large. If we use the full power of the context-free languages we get compilers which in general are inefficient, and probably not so good in handling erroneous input. This latter fact may not be so important from a theoretical point of view, but it is from a pragmatical point of view. Most invocations of compilers still have as their primary goal to discover mistakes made when typing the program, and not so much generating actual code. This aspect is even stronger present in strongly typed languages, such as Java and Haskell, where the type checking performed by the compilers is one of the main contributions to the increase in efficiency in the programming process.

When constructing a recogniser for a language described by a context-free grammar one often wants to check whether or not the grammar has specific desirable properties. Unfortunately, for a human being it is not always easy, and quite often practically impossible, to see whether or not a particular property holds. Furthermore, it may be very expensive to check whether or not such a property holds. This has led to a

1. Goals

whole hierarchy of context-free grammars classes, some of which are more powerful, some are easy to check by machine, and some are easily checked by a simple human inspection. In this course we will see many examples of such classes. The general observation is that the more precise the answer to a specific question one wants to have, the more computational effort is needed and the sooner this question cannot be answered by a human being anymore.

1.3. Compositionality

As we will see the structure of many compilers follows directly from the grammar that describes the language to be compiled. Once this phenomenon was recognised it went under the name *syntax directed compilation*. Under closer scrutiny, and under the influence of the more functional oriented style of programming, it was recognised that actually compilers are a special form of homomorphisms, a concept thus far only familiar to mathematicians and more theoretically oriented computer scientist that study the description of the meaning of a programming language.

This should not come as a surprise since this recognition is a direct consequence of the tendency that ever greater parts of compilers are more or less automatically generated from a formal description of some aspect of a programming language; e.g. by making use of a description of their outer appearance or by making use of a description of the semantics (meaning) of a language. We will see many examples of such mappings. As a side effect you will acquire a special form of writing functional programs, which makes it often surprisingly simple to solve at first sight rather complicated programming assignments. We will see that the concept of *lazy evaluation* plays an important rôle in making these efficient and straightforward implementations possible.

1.4. Abstraction mechanisms

One of the main reasons for that what used to be an endeavour for a large team in the past can now easily be done by a couple of first year's students in a matter of days or weeks, is that over the last thirty years we have discovered the right kind of abstractions to be used, and an efficient way of partitioning a problem into smaller components. Unfortunately there is no simple way to teach the techniques which have led us thus far. The only way we see is to take a historians view and to compare the old and the new situations.

Fortunately however there have also been some developments in programming language design, of which we want to mention the developments in the area of functional programming in particular. We claim that the combination of a modern, albeit quite elaborate, type system, combined with the concept of lazy evaluation, provides an ideal platform to develop and practice ones abstraction skills. There does not exist

another readily executable formalism which may serve as an equally powerful tool. We hope that by presenting many algorithms, and fragments thereof, in a modern functional language, we can show the real power of abstraction, and even find some inspiration for further developments in language design: i.e., find clues about how to extend such languages to enable us to make common patterns, which thus far have only been demonstrated by giving examples, explicit.

1. Goals

2. Context-Free Grammars

Introduction

We often want to recognise a particular structure hidden in a sequence of symbols. For example, when reading this sentence, you automatically structure it by means of your understanding of the English language. Of course, not any sequence of symbols is an English sentence. So how do we characterise English sentences? This is an old question, which was posed long before computers were widely used; in the area of natural language research the question has often been posed what actually constitutes a “language”. The simplest definition one can come up with is to say that the English language equals the set of all grammatically correct English sentences, and that a sentence consists of a sequence of English words. This terminology has been carried over to computer science: the programming language Java can be seen as the set of all correct Java programs, whereas a Java program can be seen as a sequence of Java symbols, such as identifiers, reserved words, specific operators etc.

This chapter introduces the most important notions of this course: the concept of a *language* and a *grammar*. A language is a, possibly infinite, set of sentences and sentences are sequences of symbols taken from a finite set (e.g., sequences of characters, which are referred to as strings). Just as we say that the fact whether or not a sentence belongs to the English language is determined by the English grammar (remember that before we have used the phrase “grammatically correct”), we have a grammatical formalism for describing artificial languages.

A difference with the grammars for natural languages is that this grammatical formalism is a completely formal one. This property may enable us to mathematically prove that a sentence belongs to some language, and often such proofs can be constructed automatically by a computer in a process called *parsing*. Notice that this is quite different from the grammars for natural languages, where one may easily disagree about whether something is correct English or not. This completely formal approach however also comes with a disadvantage; the expressiveness of the class of grammars we are going to describe in this chapter is rather limited, and there are many languages one might want to describe but which cannot be described, given the limitations of the formalism.

2. Context-Free Grammars

Goals

The main goal of this chapter is to introduce and show the relation between the main concepts for describing the parsing problem: languages and sentences, and grammars.

In particular, after you have studied this chapter you will:

- know the concepts of *language* and *sentence*;
- know how to describe languages by means of *context-free grammars*;
- know the difference between a *terminal* symbol and a *nonterminal* symbol;
- be able to read and interpret the *BNF* notation;
- understand the *derivation* process used in describing languages;
- understand the rôle of *parse trees*;
- understand the relation between context-free grammars and datatypes;
- understand the *EBNF* formalism;
- understand the concepts of *concrete* and *abstract syntax*;
- be able to convert a grammar from EBNF-notation into BNF-notation by hand;
- be able to construct a simple context-free grammar in EBNF notation;
- be able to verify whether or not a simple grammar is *ambiguous*;
- be able to *transform* a grammar, for example for removing left recursion.

2.1. Languages

The goal of this section is to introduce the concepts of language and sentence.

In conventional texts about mathematics it is not uncommon to encounter a definition of sequences that looks as follows:

sequence

Definition 2.1 (Sequence). Let X be a set. The set of sequences over X , called X^* , is defined as follows:

- ε is a sequence, called the empty sequence, and
- if z is a sequence and a is an element of X , then az is also a sequence.

induction

The above definition is an instance of a very common definition pattern: it is a *definition by induction*, i. e., the definition of the concept refers to the concept itself. It is implicitly understood that nothing that cannot be formed by repeated, but finite application of one of the two given rules is a sequence over X .

Furthermore, the definition corresponds almost exactly to the definition of the type $[a]$ of *lists* with elements of type a in Haskell. The one difference is that Haskell lists can be infinite, whereas sequences are always finite.

In the following, we will introduce several concepts based on sequences. They can be implemented easily in Haskell using lists.

Functions that operate on an inductively defined structure such as sequences are typically *structurally recursive*, i. e., such definitions often follow a recursion pattern which is similar to the definition of the structure itself. (Recall the function *foldr* from the course on Functional Programming.)

Note that the Haskell notation for lists is generally more precise than the mathematical notation for sequences. When talking about languages and grammars, we often leave the distinction between single symbols and sequences implicit.

We use letters from the beginning of the alphabet to represent single symbols, and letters from the end of the alphabet to represent sequences. We write a to denote both the single symbol a or the sequence $a\varepsilon$, depending on context. We typically use ε only when we want to explicitly emphasize that we are talking about the empty sequence.

Furthermore, we denote concatenation of sequences and symbols in the same way, i. e., az should be understood as the symbol a followed by the sequence z , whereas xy is the concatenation of sequences x and y .

In Haskell, all these distinctions are explicit. Elements are distinguished from lists by their type; there is a clear difference between a and $[a]$. Concatenation of lists is handled by the operator $(++)$, whereas a single element can be added to the front of a list using $(:)$. Also, Haskell identifiers often have longer names, so ab in Haskell is to be understood as a single identifier with name ab , not as a combination of two symbols a and b .

Now we move from individual sequences to finite or infinite sets of sequences. We start with some terminology:

Definition 2.2 (Alphabet, Language, Sentence).

- An *alphabet* is a finite set of *symbols*.
- A *language* is a subset of T^* , for some alphabet T .
- A *sentence* (often also called *word*) is an element of a language.

alphabet
language
sentence

Note that ‘word’ and ‘sentence’ in formal languages are used as synonyms.

Some examples of alphabets are:

- the conventional Roman alphabet: $\{\mathbf{a, b, c, \dots, z}\}$;
- the binary alphabet $\{0, 1\}$;
- sets of reserved words $\{\mathbf{if, then, else}\}$;
- a set of characters $l = \{\mathbf{a, b, c, d, e, i, k, l, m, n, o, p, r, s, t, u, w, x}\}$;
- a set of English words $\{\mathbf{course, practical, exercise, exam}\}$.

Examples of languages are:

- T^* , \emptyset (the empty set), $\{\varepsilon\}$ and T are languages over alphabet T ;

2. Context-Free Grammars

- the set $\{\text{course, practical, exercise, exam}\}$ is a language over the alphabet l of characters and exam is a sentence in it.

The question that now arises is how to *specify* a language. Since a language is a set we immediately see three different approaches:

- enumerate all the elements of the set explicitly;
- characterise the elements of the set by means of a predicate;
- define which elements belong to the set by means of induction.

We have just seen some examples of the first (the Roman alphabet) and third (the set of sequences over an alphabet) approach. Examples of the second approach are:

- the even natural numbers $\{n \mid n \in \{0, 1, \dots, 9\}^*, n \bmod 2 = 0\}$;
- the language PAL of palindromes, sequences which read the same forward as backward, over the alphabet $\{a, b, c\}$: $\{s \mid s \in \{a, b, c\}^*, s = s^R\}$, where s^R denotes the reverse of sequence s .

One of the fundamental differences between the predicative and the inductive approach to defining a language is that the latter approach is *constructive*, i.e., it provides us with a way to enumerate all elements of a language. If we define a language by means of a predicate we only have a means to decide whether or not an element belongs to a language. A famous example of a language which is easily defined in a predicative way, but for which the membership test is very hard, is the set of prime numbers.

Quite often we want to prove that a language L , which is defined by means of an inductive definition, has a specific property P . If this property is of the form $P(L) = \forall x \in L. P(x)$, then we want to prove that $L \subseteq P$.

Since languages are sets the usual set operators such as union, intersection and difference can be used to construct new languages from existing ones. The complement of a language L over alphabet T is defined by $\bar{L} = \{x \mid x \in T^*, x \notin L\}$.

In addition to these set operators, there are more specific operators, which apply only to sets of sequences. We will use these operators mainly in the chapter on regular languages, Chapter 8. Note that \cup denotes set union, so $\{1, 2\} \cup \{1, 3\} = \{1, 2, 3\}$.

Definition 2.3 (Language operations). Let L and M be languages over the same alphabet T , then

\bar{L}	$= T^* - L$	complement of L
L^R	$= \{s^R \mid s \in L\}$	reverse of L
LM	$= \{st \mid s \in L, t \in M\}$	concatenation of L and M
L^0	$= \{\varepsilon\}$	0 th power of L
L^{n+1}	$= LL^n$	$n + 1^{\text{st}}$ power of L
L^*	$= \bigcup_{i \in \mathbb{N}} L^i = L^0 \cup L^1 \cup L^2 \cup \dots$	star-closure of L
L^+	$= \bigcup_{i \in \mathbb{N}, i > 0} L^i = L^1 \cup L^2 \cup \dots$	positive closure of L

The following equations follow immediately from the above definitions.

$$\begin{aligned} L^* &= \{\varepsilon\} \cup LL^* \\ L^+ &= LL^* \end{aligned}$$

Exercise 2.1. Let $L = \{\text{ab, aa, baa}\}$, where **a** and **b** are the terminals. Which of the following strings are in L^* : **abaabaaabaa**, **aaaabaaaa**, **baaaaabaaaab**, **baaaaabaa**?

Exercise 2.2. What are the elements of \emptyset^* ?

Exercise 2.3. For any language, prove

1. $\emptyset L = L\emptyset = \emptyset$
2. $\{\varepsilon\}L = L\{\varepsilon\} = L$

Exercise 2.4. In this section we defined two “star” operators: one for arbitrary sets (Definition 2.1) and one for languages (Definition 2.3). Is there a difference between these operators?

2.2. Grammars

The goal of this section is to introduce the concept of context-free grammars.

Working with sets might be fun, but it often is complicated to manipulate sets, and to prove properties of sets. For these purposes we introduce syntactical definitions, called grammars, of sets. This section will only discuss so-called *context-free grammars*, a kind of grammars that are convenient for automatic processing, and that can describe a large class of languages. But the class of languages that can be described by context-free grammars is limited.

In the previous section we defined PAL, the language of palindromes, by means of a predicate. Although this definition defines the language we want, it is hard to use in proofs and programs. An important observation is the fact that the set of palindromes can be defined inductively as follows.

Definition 2.4 (Palindromes by induction).

- The empty string, ε , is a palindrome;
- the strings consisting of just one character, **a**, **b**, and **c**, are palindromes;
- if P is a palindrome, then the strings obtained by prepending and appending the same character, **a**, **b**, and **c**, to it are also palindromes, that is, the strings

$$\begin{aligned} &\mathbf{aPa} \\ &\mathbf{bPb} \\ &\mathbf{cPc} \end{aligned}$$

are palindromes.

2. Context-Free Grammars

The first two parts of the definition cover the basic cases. The last part of the definition covers the inductive cases. All strings which belong to the language PAL inductively defined using the above definition read the same forwards and backwards. Therefore this definition is said to be *sound* (every string in PAL is a palindrome). Conversely, if a string consisting of a's, b's, and c's reads the same forwards and backwards then it belongs to the language PAL. Therefore this definition is said to be *complete* (every palindrome is in PAL).

Finding an inductive definition for a language which is described by a predicate (like the one for palindromes) is often a nontrivial task. Very often it is relatively easy to find a definition that is sound, but you also have to convince yourself that the definition is complete. A typical method for proving soundness and completeness of an inductive definition is *mathematical induction*.

Now that we have an inductive definition for palindromes, we can proceed by giving a formal representation of this inductive definition.

Inductive definitions like the one above can be represented formally by making use of *deduction rules* which look like:

$$a_1, a_2, \dots, a_n \vdash a \quad \text{or} \quad \vdash a$$

The first kind of deduction rule has to be read as follows:

if a_1, a_2, \dots and a_n are true,
then a is true.

The second kind of deduction rule, called an axiom, has to be read as follows:

a is true.

Using these deduction rules we can now write down the inductive definition for PAL as follows:

$$\begin{aligned} &\vdash \varepsilon \in \text{PAL}' \\ &\vdash \mathbf{a} \in \text{PAL}' \\ &\vdash \mathbf{b} \in \text{PAL}' \\ &\vdash \mathbf{c} \in \text{PAL}' \\ &P \in \text{PAL}' \vdash \mathbf{aPa} \in \text{PAL}' \\ &P \in \text{PAL}' \vdash \mathbf{bPb} \in \text{PAL}' \\ &P \in \text{PAL}' \vdash \mathbf{cPc} \in \text{PAL}' \end{aligned}$$

Although the definition of PAL' is completely formal, it is still laborious to write. Since in computer science we use many definitions which follow such a pattern, we introduce a shorthand for it, called a *grammar*. A grammar consists of production rules. We can give a grammar of PAL' by translating the deduction rules given above into production rules. The rule with which the empty string is constructed is:

$$P \rightarrow \varepsilon$$

This rule corresponds to the axiom that states that the empty string ε is a palindrome. A rule of the form $s \rightarrow \alpha$, where s is symbol and α is a sequence of symbols, is called a *production rule*, or *production* for short. A production rule can be considered as a possible way to rewrite the symbol s . The symbol P to the left of the arrow is a symbol which denotes palindromes. Such a symbol is an example of a *nonterminal symbol*, or *nonterminal* for short. Nonterminal symbols are also called auxiliary symbols: their only purpose is to denote structure, they are not part of the alphabet of the language. Three other basic production rules are the rules for constructing palindromes consisting of just one character. Each of the one element strings a , b , and c is a palindrome, and gives rise to a production:

$$P \rightarrow a$$

$$P \rightarrow b$$

$$P \rightarrow c$$

These production rules correspond to the axioms that state that the one element strings a , b , and c are palindromes. If a string α is a palindrome, then we obtain a new palindrome by prepending and appending an a , b , or c to it, that is, $a\alpha a$, $b\alpha b$, and $c\alpha c$ are also palindromes. To obtain these palindromes we use the following recursive productions:

$$P \rightarrow aPa$$

$$P \rightarrow bPb$$

$$P \rightarrow cPc$$

These production rules correspond to the deduction rules that state that, if P is a palindrome, then one can deduce that aPa , bPb and cPc are also palindromes. The grammar we have presented so far consists of three components:

- the set of *terminals* $\{a, b, c\}$;
- the set of *nonterminals* $\{P\}$;
- and the set of productions (the seven productions that we have introduced so far).

Note that the intersection of the set of terminals and the set of nonterminals is empty. We complete the description of the grammar by adding a fourth component: the nonterminal *start symbol* P . In this case we have only one choice for a start symbol, but a grammar may have many nonterminal symbols, and we always have to select one to start with.

To summarize, we obtain the following grammar for PAL:

$$P \rightarrow \varepsilon$$

$$P \rightarrow a$$

$$P \rightarrow b$$

2. Context-Free Grammars

$$\begin{aligned}P &\rightarrow c \\P &\rightarrow aPa \\P &\rightarrow bPb \\P &\rightarrow cPc\end{aligned}$$

The definition of the set of terminals, $\{a, b, c\}$, and the set of nonterminals, $\{P\}$, is often implicit. Also the start-symbol is implicitly defined here since there is only one nonterminal.

We conclude this example with the formal definition of a context-free grammar.

context-free
grammar

Definition 2.5 (Context-Free Grammar). A context-free grammar G is a four-tuple (T, N, R, S) where

- T is a finite set of terminal symbols;
- N is a finite set of nonterminal symbols (T and N are disjoint);
- R is a finite set of production rules. Each production has the form $A \rightarrow \alpha$, where A is a nonterminal and α is a sequence of terminals and nonterminals;
- S is the start symbol, $S \in N$.

The adjective “context-free” in the above definition comes from the specific production rules that are considered: exactly one nonterminal on the left hand side. Not every language can be described via a context-free grammar. The standard example here is $\{a^n b^n c^n \mid n \in \mathbb{N}\}$. We will encounter this example again later in these lecture notes.

2.2.1. Notational conventions

In the definition of the grammar for PAL we have written every production on a single line. Since this takes up a lot of space, and since the production rules form the heart of every grammar, we introduce the following shorthand. Instead of writing

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

we combine the two productions for S in one line as using the symbol $|$:

$$S \rightarrow \alpha \mid \beta$$

We may rewrite any number of rewrite rules for one nonterminal in this fashion, so the grammar for PAL may also be written as follows:

$$P \rightarrow \varepsilon \mid a \mid b \mid c \mid aPa \mid bPb \mid cPc$$

BNF

The notation we use for grammars is known as *BNF* – Backus Naur Form – after Backus and Naur, who first used this notation for defining grammars.

Another notational convention concerns names of productions. Sometimes we want to give names to production rules. The names will be written in front of the production. So, for example,

$$\begin{aligned} \text{Alpha rule: } & S \rightarrow \alpha \\ \text{Beta rule: } & S \rightarrow \beta \end{aligned}$$

Finally, if we give a context-free grammar just by means of its productions, the start-symbol is usually the nonterminal in the left hand side of the first production, and the start-symbol is usually called S .

Exercise 2.5. Give a context free grammar for the set of sentences over alphabet X where

1. $X = \{\mathbf{a}\}$,
2. $X = \{\mathbf{a}, \mathbf{b}\}$.

Exercise 2.6. Give a context free grammar for the language

$$L = \{a^n b^n \mid n \in \mathbb{N}\}$$

Exercise 2.7. Give a grammar for *palindromes* over the alphabet $\{\mathbf{a}, \mathbf{b}\}$.

Exercise 2.8. Give a grammar for the language

$$L = \{s s^R \mid s \in \{\mathbf{a}, \mathbf{b}\}^*\}$$

This language is known as the language of *mirror palindromes*.

Exercise 2.9. A *parity sequence* is a sequence consisting of 0's and 1's that has an even number of ones. Give a grammar for parity sequences.

Exercise 2.10. Give a grammar for the language

$$L = \{w \mid w \in \{\mathbf{a}, \mathbf{b}\}^* \wedge \#(\mathbf{a}, w) = \#(\mathbf{b}, w)\}$$

where $\#(c, w)$ is the number of c -occurrences in w .

2.3. The language of a grammar

The goal of this section is to describe the relation between grammars and languages: to show how to derive sentences of a language, given its grammar.

In the previous section, we have demonstrated in several examples how to construct a grammar for a particular language. Now we consider the reverse question: how to obtain a language from a given grammar? Before we can answer this question we first have to say what we can do with a grammar. The answer is simple: we can derive sequences with it.

2. Context-Free Grammars

How do we construct a palindrome? A palindrome is a sequence of terminals, in our case the characters **a**, **b** and **c**, that can be derived in zero or more direct derivation steps from the start symbol P using the productions of the grammar for palindromes given before.

For example, the sequence **bacab** can be derived using the grammar for palindromes as follows:

$$\begin{aligned} & P \\ \Rightarrow & \mathbf{bPb} \\ \Rightarrow & \mathbf{baPab} \\ \Rightarrow & \mathbf{bacab} \end{aligned}$$

derivation

Such a construction is called a *derivation*. In the first step of this derivation production $P \rightarrow \mathbf{bPb}$ is used to rewrite P into \mathbf{bPb} . In the second step production $P \rightarrow \mathbf{aPa}$ is used to rewrite \mathbf{bPb} into \mathbf{baPab} . Finally, in the last step production $P \rightarrow \mathbf{c}$ is used to rewrite \mathbf{baPab} into \mathbf{bacab} . Constructing a derivation can be seen as a constructive proof that the string **bacab** is a palindrome.

We will now describe derivation steps more formally.

Definition 2.6 (Derivation). Suppose $X \rightarrow \beta$ is a production of a grammar, where X is a nonterminal symbol and β is a sequence of (nonterminal or terminal) symbols. Let $\alpha X \gamma$ be a sequence of (nonterminal or terminal) symbols. We say that $\alpha X \gamma$ *directly derives* the sequence $\alpha \beta \gamma$, which is obtained by replacing the left hand side X of the production by the corresponding right hand side β . We write $\alpha X \gamma \Rightarrow \alpha \beta \gamma$ and we also say that $\alpha X \gamma$ rewrites to $\alpha \beta \gamma$ in one step. A sequence φ_n is *derived* from a sequence φ_0 , written $\varphi_0 \Rightarrow^* \varphi_n$, if there exist sequences $\varphi_0, \dots, \varphi_n$ such that

direct derivation

derivation

$$\forall i, 0 \leq i < n : \quad \varphi_i \Rightarrow \varphi_{i+1}$$

If $n = 0$, this statement is trivially true, and it follows that we can derive each sentence φ from itself in zero steps:

$$\varphi \Rightarrow^* \varphi$$

partial derivation

A *partial derivation* is a derivation of a sequence that still contains nonterminals.

Finding a derivation $\varphi_0 \Rightarrow^* \varphi_n$ is, in general, a nontrivial task. A derivation is only one branch of a whole search tree which contains many more branches. Each branch represents a (successful or unsuccessful) direction in which a possible derivation may proceed. Another important challenge is to arrange things in such a way that finding a derivation can be done in an efficient way.

From the example derivation above it follows that

$$P \Rightarrow^* \text{bacab}$$

Because this derivation begins with the start symbol of the grammar and results in a sequence consisting of terminals only (a terminal string), we say that the string **bacab** belongs to the language generated by the grammar for palindromes. In general, we define

Definition 2.7 (Language of a grammar). The *language of a grammar* $G=(T, N, R, S)$, usually denoted by $L(G)$, is defined as language of a grammar

$$L(G) = \{s \mid S \Rightarrow^* s, s \in T^*\}$$

The language $L(G)$ is also called *the language generated by the grammar* G . We sometimes also talk about the language of a nonterminal A , which is defined by

$$L(A) = \{s \mid A \Rightarrow^* s, s \in T^*\}$$

Note that different grammars may have the same language. For example, if we extend the grammar for PAL with the production $P \rightarrow \text{bacab}$, we obtain a grammar with exactly the same language as PAL. Two grammars that generate the same language are called *equivalent*. So for a particular grammar there exists a unique language, but the reverse is not true: given a language we can construct many grammars that generate the language. To phrase it more mathematically: the mapping between a grammar and its language is not a bijection. equivalent

Definition 2.8 (Context-free language). A *context-free language* is a language that is generated by a context-free grammar. context-free language

All palindromes can be derived from the start symbol P . Thus, the language of our grammar for palindromes is PAL, the set of all palindromes over the alphabet $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$, and PAL is context-free.

2.3.1. Examples of basic languages

Digits occur in a several programming languages and other languages, and so do letters. In this subsection we will define some grammars that specify some basic languages such as digits and letters. These grammars will be used frequently in later sections.

- The language of single digits is specified by a grammar with ten production rules for the nonterminal *Dig*.

$$\text{Dig} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

2. Context-Free Grammars

- We obtain sequences of digits by means of the following grammar:

$$Digs \rightarrow \varepsilon \mid Dig\ Digs$$

- Natural numbers are sequences of digits that start with a non-zero digit. So in order to specify natural numbers, we first define the language of non-zero digits.

$$Dig-0 \rightarrow 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Now we can define the language of natural numbers as follows.

$$Nat \rightarrow 0 \mid Dig-0\ Digs$$

- Integers are natural numbers preceded by a sign. If a natural number is not preceded by a sign, it is supposed to be a positive number.

$$\begin{aligned} Sign &\rightarrow + \mid - \\ Int &\rightarrow Sign\ Nat \mid Nat \end{aligned}$$

- The languages of small letters and capital letters are each specified by a grammar with 26 productions:

$$\begin{aligned} SLetter &\rightarrow a \mid b \mid \dots \mid z \\ CLetter &\rightarrow A \mid B \mid \dots \mid Z \end{aligned}$$

In the real definitions of these grammars we have to write each of the 26 letters, of course. A letter is now either a small or a capital letter.

$$Letter \rightarrow SLetter \mid CLetter$$

- Variable names, function names, datatypes, etc., are all represented by identifiers in programming languages. The following grammar for identifiers might be used in a programming language:

$$\begin{aligned} Identifier &\rightarrow Letter\ AlphaNums \\ AlphaNums &\rightarrow \varepsilon \mid Letter\ AlphaNums \mid Dig\ AlphaNums \end{aligned}$$

An identifier starts with a letter, and is followed by a sequence of alphanumeric characters, i. e., letters and digits. We might want to allow more symbols, such as for example underscores and dollar symbols, but then we have to adjust the grammar, of course.

- Dutch zip codes consist of four digits, of which the first digit is non-zero, followed by two capital letters. So

$$ZipCode \rightarrow Dig-0\ Dig\ Dig\ Dig\ CLetter\ CLetter$$

Exercise 2.11. A terminal string that belongs to the language of a grammar is always derived in one or more steps from the start symbol of the grammar. Why?

Exercise 2.12. What language is generated by the grammar with the single production rule

$$S \rightarrow \varepsilon$$

Exercise 2.13. What language does the grammar with the following productions generate?

$$\begin{aligned} S &\rightarrow Aa \\ A &\rightarrow B \\ B &\rightarrow Aa \end{aligned}$$

Exercise 2.14. Give a simple description of the language generated by the grammar with productions

$$\begin{aligned} S &\rightarrow aA \\ A &\rightarrow bS \\ S &\rightarrow \varepsilon \end{aligned}$$

Exercise 2.15. Is the language L defined in Exercise 2.1 context free ?

2.4. Parse trees

The goal of this section is to introduce parse trees, and to show how parse trees relate to derivations. Furthermore, this section defines (non)ambiguous grammars.

For any partial derivation, i. e., a derivation that contains nonterminals in its right hand side, there may be several productions of the grammar that can be used to proceed the partial derivation with. As a consequence, there may be different derivations for the same sentence.

However, if only the order in which the derivation steps are chosen differs between two derivations, then the derivations are considered to be equivalent. If, however, different derivation steps have been chosen in two derivations, then these derivations are considered to be different.

Here is a simple example. Consider the grammar SequenceOfS with productions:

$$\begin{aligned} S &\rightarrow SS \\ S &\rightarrow s \end{aligned}$$

Using this grammar, we can derive the sentence sss in at least the following two ways (the nonterminal that is rewritten in each step appears underlined):

$$\underline{S} \Rightarrow \underline{S}\underline{S} \Rightarrow \underline{S}\underline{S}\underline{S} \Rightarrow \underline{S}s\underline{S} \Rightarrow \underline{ss}\underline{S} \Rightarrow \underline{sss}$$

2. Context-Free Grammars

$$\underline{S} \Rightarrow \underline{SS} \Rightarrow \underline{sS} \Rightarrow \underline{sSS} \Rightarrow \underline{sSs} \Rightarrow \underline{sss}$$

These derivations are the same up to the order in which derivation steps are taken. However, the following derivation does not use the same derivation steps:

$$\underline{S} \Rightarrow \underline{SS} \Rightarrow \underline{SSS} \Rightarrow \underline{sSS} \Rightarrow \underline{ssS} \Rightarrow \underline{sss}$$

In both derivation sequences above, the first S was rewritten to s . In this derivation, however, the first S is rewritten to SS .

leftmost derivation

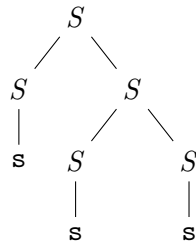
The set of all equivalent derivations can be represented by selecting a so-called *canonical element*. A good candidate for such a canonical element is the *leftmost derivation*. In a leftmost derivation, the leftmost nonterminal is rewritten in each step. If there exists a derivation of a sentence x using the productions of a grammar, then there exists also a leftmost derivation of x . The last of the three derivation sequences for the sentence sss given above is a leftmost derivation. The two equivalent derivation sequences before, however, are both not leftmost. The leftmost derivation corresponding to these two sequences above is

$$\underline{S} \Rightarrow \underline{SS} \Rightarrow \underline{sS} \Rightarrow \underline{sSS} \Rightarrow \underline{ssS} \Rightarrow \underline{sss}$$

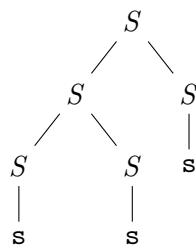
parse tree

There exists another convenient way for representing equivalent derivations: they all correspond to the same *parse tree* (or derivation tree). A parse tree is a representation of a derivation which abstracts from the order in which derivation steps are chosen. The internal nodes of a parse tree are labelled with a nonterminal N , and the children of such a node are the parse trees for symbols of the right hand side of a production for N . The parse tree of a terminal symbol is a leaf labelled with the terminal symbol.

The resulting parse tree of the first two derivations of the sentence sss looks as follows:



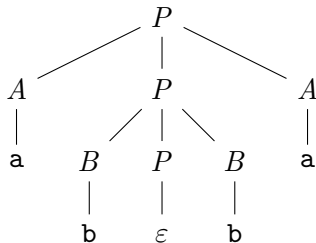
The third derivation of the sentence sss results in a different parse tree:



As another example, all derivations of the string **abba** using the productions of the grammar

$$\begin{aligned} P &\rightarrow \varepsilon \\ P &\rightarrow APA \\ P &\rightarrow BPB \\ A &\rightarrow \mathbf{a} \\ B &\rightarrow \mathbf{b} \end{aligned}$$

are represented by the following derivation tree:



A derivation tree can be seen as a *structural interpretation* of the derived sentence. Note that there might be more than one structural interpretation of a sentence with respect to a given grammar. Such grammars are called *ambiguous*.

Definition 2.9 (ambiguous grammar, unambiguous grammar). A grammar is *unambiguous* if every sentence has a unique leftmost derivation, or, equivalently, if every sentence has a unique derivation tree. Otherwise it is called *ambiguous*.

unambiguous
ambiguous

The grammar SequenceOfS for constructing sequences of **s**'s is an example of an ambiguous grammar, since there exist two parse trees for the sentence **sss**.

It is in general undecidable whether or not an arbitrary context-free grammar is ambiguous. This implies that it is impossible to write a program that determines for an arbitrary context-free grammar whether it is ambiguous or not.

It is usually rather difficult to translate languages with ambiguous grammars. Therefore, you will find that most grammars of programming languages and other languages that are used in processing information are unambiguous.

Grammars have proved very successful in the specification of artificial languages (such as programming languages). They have proved less successful in the specification of natural languages (such as English), partly because it is extremely difficult to construct an unambiguous grammar that specifies a nontrivial part of the language. Take for example the sentence 'They are flying planes'. This sentence can be read in two ways, with different meanings: 'They – are – flying planes', and 'They – are flying – planes'. While ambiguity of natural languages may perhaps be considered as an advantage for their users (e. g., politicians), it certainly is considered a disadvantage for language translators, because it is usually impossible to maintain an ambiguous meaning in a translation.

2.5. Grammar transformations

In this section, we will first look at a number of properties of grammars. We then will show how we can systematically transform grammars into grammars that describe the same language and satisfy particular properties.

Here are some examples of properties that are worth considering for a particular grammar:

- a grammar may be unambiguous, that is, every sentence of its language has a unique parse tree;
- a grammar may have the property that only the start symbol can derive the empty string; no other nonterminal can derive the empty string;
- a grammar may have the property that every production either has a single terminal, or two nonterminals in its right hand side. Such a grammar is said to be in Chomsky normal form.

So why are we interested in such properties? Some of these properties imply that it is possible to build parse trees for sentences of the language of the grammar in only one way. Some other properties imply that we can build these parse trees very fast. Other properties are used to prove facts about grammars. Yet other properties are used to efficiently compute certain information from parse trees of a grammar.

Properties are particularly interesting in combination with grammar transformations. A *grammar transformation* is a procedure to obtain a grammar G' from a grammar G such that $L(G') = L(G)$.

grammar
transformation

Suppose, for example, that we have a program that builds parse trees for sentences of grammars in Chomsky normal form, and that we can prove that every grammar can be transformed in a grammar in Chomsky normal form. Then we can use this program for building parse trees for any grammar.

Since it is sometimes convenient to have a grammar that satisfies a particular property for a language, we would like to be able to transform grammars into other grammars that generate the same language, but that possibly satisfy different properties. In the following, we describe a number of grammar transformations:

- Removing duplicate productions.
- Substituting right hand sides for nonterminals.
- Removing unreachable productions.
- Left factoring.
- Removing left recursion.
- Associative separator.
- Introduction of priorities.

There are many more transformations than we describe here; we will only show a small but useful set of grammar transformations. In the following, we will assume that N is the set of nonterminals, T is the set of terminals, and that u, v, w, x, y and z denote sequences of terminals and nonterminals, i. e., are elements of $(N \cup T)^*$.

2.5.1. Removing duplicate productions

This grammar transformation is a transformation that can be applied to any grammar of the correct form. If a grammar contains two occurrences of the same production rule, one of these occurrences can be removed. For example,

$$A \rightarrow u \mid u \mid v$$

can be transformed into

$$A \rightarrow u \mid v$$

2.5.2. Substituting right hand sides for nonterminals

If a nonterminal X occurs in a right hand side of a production, the production may be replaced by just as many productions as there exist productions for X , in which X has been replaced by its right hand sides. For example, we can substitute B in the right hand side of the first production in the following grammar:

$$\begin{aligned} A &\rightarrow uBv \mid z \\ B &\rightarrow x \mid w \end{aligned}$$

The resulting grammar is:

$$\begin{aligned} A &\rightarrow uxv \mid uuv \mid z \\ B &\rightarrow x \mid w \end{aligned}$$

2.5.3. Removing unreachable productions

Consider the result of the transformation above:

$$\begin{aligned} A &\rightarrow uxv \mid uuv \mid z \\ B &\rightarrow x \mid w \end{aligned}$$

If A is the start symbol and B does not occur in any of the symbol sequences u, v, w, x, z , then the second production can never occur in the derivation of a sentence starting from A . In such a case, the unreachable production can be dropped:

$$A \rightarrow uxv \mid uuv \mid z$$

2.5.4. Left factoring

left factoring

Left factoring is a grammar transformation that is applicable when two productions for the same nonterminal start with the same sequence of (terminal and/or nonterminal) symbols. These two productions can then be replaced by a single production, that ends with a new nonterminal, replacing the part of the sequence after the common start sequence. Two productions for the new nonterminal are added: one for each of the two different end sequences of the two productions. For example:

$$A \rightarrow xy \mid xz \mid v$$

may be transformed into

$$\begin{aligned} A &\rightarrow xZ \mid v \\ Z &\rightarrow y \mid z \end{aligned}$$

where Z is a new nonterminal. As we will see in Chapter 3, parsers can be constructed systematically from a grammar, and left factoring the grammar before constructing a parser can dramatically improve the performance of the resulting parser.

2.5.5. Removing left recursion

left recursion

A production is called *left-recursive* if the right-hand side starts with the nonterminal of the left-hand side. For example, the production

$$A \rightarrow Az$$

is left-recursive. A grammar is left-recursive if we can derive $A \Rightarrow^+ Az$ for some nonterminal A of the grammar (i.e., if we can derive Az from A in one or more steps).

Left-recursive grammars are sometimes undesirable – we will, for instance, see in Chapter 3 that a parser constructed systematically from a left-recursive grammar may loop. Fortunately, left recursion can be removed by transforming the grammar. The following transformation removes left-recursive productions.

To remove the left-recursive productions of a nonterminal A , we divide the productions for A in sets of left-recursive and non left-recursive productions. We can thus factorize the productions for A as follows:

$$\begin{aligned} A &\rightarrow Ax_1 \mid Ax_2 \mid \dots \mid Ax_n \\ A &\rightarrow y_1 \mid y_2 \mid \dots \mid y_m \end{aligned}$$

where none of the symbol sequences y_1, \dots, y_m starts with A . We now add a new nonterminal Z , and replace the productions for A by:

$$\begin{aligned} A &\rightarrow y_1 \mid y_1 Z \mid \dots \mid y_m \mid y_m Z \\ Z &\rightarrow x_1 \mid x_1 Z \mid \dots \mid x_n \mid x_n Z \end{aligned}$$

Note that this procedure only works for a grammar that is *directly* left-recursive, i. e., a grammar that contains a left-recursive production of the form $A \rightarrow Ax$.

Grammars can also be indirectly left recursive. An example is

$$\begin{aligned} A &\rightarrow Bx \\ B &\rightarrow Ay \end{aligned}$$

None of the two productions is left recursive, but we can still derive $A \Rightarrow^* Ayx$. Removing left recursion in an indirectly left-recursive grammar is also possible, but a bit more complicated [1].

Here is an example of how to apply the procedure described above to a grammar that is directly left recursive, namely the grammar `SequenceOfS` that we have introduced in Section 2.4:

$$\begin{aligned} S &\rightarrow SS \\ S &\rightarrow s \end{aligned}$$

The first production is left-recursive. The second is not. We can thus directly apply the procedure for left recursion removal, and obtain the following productions:

$$\begin{aligned} S &\rightarrow s \mid sZ \\ Z &\rightarrow S \mid SZ \end{aligned}$$

2.5.6. Associative separator

The following grammar fragment generates a list of declarations, separated by a semicolon ‘;’:

$$\begin{aligned} Decls &\rightarrow Decls ; Decls \\ Decls &\rightarrow Decl \end{aligned}$$

The productions for *Decl*, which generates a single declaration, have been omitted. This grammar is ambiguous, for the same reason as `SequenceOfS` is ambiguous. The operator `;` is an associative separator in the generated language, that is, it does not matter how we group the declarations; given three declarations d_1 , d_2 , and d_3 , the meaning of $d_1 ; (d_2 ; d_3)$ and $(d_1 ; d_2) ; d_3$ is the same. Therefore, we may use the following unambiguous grammar for generating a language of declarations:

$$\begin{aligned} Decls &\rightarrow Decl ; Decls \\ Decls &\rightarrow Decl \end{aligned}$$

2. Context-Free Grammars

Note that in this case, the transformed grammar is also no longer left-recursive.

An alternative unambiguous (but still left-recursive) grammar for the same language is

$$\begin{aligned}Decls &\rightarrow Decls ; Decl \\Decls &\rightarrow Decl\end{aligned}$$

The grammar transformation just described should be handled with care: if the separator is associative in the generated language, like the semicolon in this case, applying the transformation is fine. However, if the separator is not associative, then removing the ambiguity in favour of a particular nesting is dangerous.

This grammar transformation is often useful for expressions which are separated by associative operators, such as for example natural numbers and addition.

2.5.7. Introduction of priorities

Another form of ambiguity often arises in the part of a grammar for a programming language which describes expressions. For example, the following grammar generates arithmetic expressions:

$$\begin{aligned}E &\rightarrow E + E \\E &\rightarrow E * E \\E &\rightarrow (E) \\E &\rightarrow Digs\end{aligned}$$

where *Digs* generates a list of digits as described in Section 2.3.1.

This grammar is ambiguous: for example, the sentence $2+4*6$ has two parse trees: one corresponding to $(2+4)*6$, and one corresponding to $2+(4*6)$. If we make the usual assumption that $*$ has higher priority than $+$, the latter expression is the intended reading of the sentence $2+4*6$. In order to obtain parse trees that respect these priorities, we transform the grammar as follows:

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

This grammar generates the same language as the previous grammar for expressions, but it respects the priorities of the operators.

In practice, often more than two levels of priority are used. Then, instead of writing a large number of nearly identically formed production rules, we can abbreviate the grammar by using parameterised nonterminals. For $1 \leq i < n$, we get productions

$$\begin{aligned} E_i &\rightarrow E_{i+1} \\ E_i &\rightarrow E_i \text{ Op}_i E_{i+1} \end{aligned}$$

The nonterminal Op_i is parameterised and generates operators of priority i . In addition to the above productions, there should also be a production for expressions of the highest priority, for example:

$$E_n \rightarrow (E_1) \mid \text{Digs}$$

2.5.8. Discussion

We have presented several examples of grammar transformations. A grammar transformation transforms a grammar into another grammar that generates the same language. For each of the above transformations we should therefore prove that the generated language remains the same. Since the proofs are too complicated at this point, they are omitted. Proofs can be found in any of the theoretical books on language and parsing theory [11].

There exist many other grammar transformations, but the ones given in this section suffice for now. Note that everywhere we use ‘left’ (left-recursion, left factoring), we can replace it by ‘right’, and obtain a dual grammar transformation. We will discuss a larger example of a grammar transformation after the following section.

Exercise 2.16. Consider the following ambiguous grammar with start symbol A :

$$\begin{aligned} A &\rightarrow A\mathbf{a}A \\ A &\rightarrow \mathbf{b} \mid \mathbf{c} \end{aligned}$$

Transform the grammar by applying the rule for associative separators. Choose the transformation such that the resulting grammar is also no longer left-recursive.

Exercise 2.17. The standard example of ambiguity in programming languages is the *dangling else*. Let G be a grammar with terminal set $\{\mathbf{if}, \mathbf{b}, \mathbf{then}, \mathbf{else}, \mathbf{a}\}$ and the following productions:

$$\begin{aligned} S &\rightarrow \mathbf{if} \ \mathbf{b} \ \mathbf{then} \ S \ \mathbf{else} \ S \\ S &\rightarrow \mathbf{if} \ \mathbf{b} \ \mathbf{then} \ S \\ S &\rightarrow \mathbf{a} \end{aligned}$$

1. Give two parse trees for the sentence $\mathbf{if} \ \mathbf{b} \ \mathbf{then} \ \mathbf{if} \ \mathbf{b} \ \mathbf{then} \ \mathbf{a} \ \mathbf{else} \ \mathbf{a}$.
2. Give an unambiguous grammar that generates the same language as G .
3. How does Java prevent this *dangling else* problem?

2. Context-Free Grammars

Exercise 2.18. A *bit list* is a nonempty list of bits separated by commas. A grammar for bit lists is given by

$$\begin{aligned}L &\rightarrow B \\L &\rightarrow L , L \\B &\rightarrow 0 \mid 1\end{aligned}$$

Remove the left recursion from this grammar.

Exercise 2.19. Consider the following grammar with start symbol S :

$$\begin{aligned}S &\rightarrow AB \\A &\rightarrow \varepsilon \mid \mathbf{a}aA \\B &\rightarrow \varepsilon \mid B\mathbf{b}\end{aligned}$$

1. What language does this grammar generate?
2. Give an equivalent non left recursive grammar.

2.6. Concrete and abstract syntax

In this section, we establish a connection between context-free grammars and Haskell datatypes. To this end, we introduce the notion of abstract syntax, and show how to obtain an abstract syntax from a concrete syntax.

For each context-free grammar we can define a corresponding datatype in Haskell. Values of these datatypes represent parse trees of the context-free grammar. As an example we take the grammar `SequenceOfS`:

$$\begin{aligned}S &\rightarrow SS \\S &\rightarrow \mathbf{s}\end{aligned}$$

First, we give each of the productions of this grammar a name:

$$\begin{aligned}\text{Beside: } S &\rightarrow SS \\ \text{Single: } S &\rightarrow \mathbf{s}\end{aligned}$$

Now we interpret the start symbol of the grammar S as a datatype, using the names of the productions as constructors:

$$\begin{aligned}\mathbf{data} S &= \text{Beside } S S \\ & \mid \text{Single}\end{aligned}$$

Note that the nonterminals on the right hand side of `Beside` reappear as arguments of the constructor `Beside`. On the other hand, the terminal symbol `s` in the production `Single` is omitted in the definition of the constructor `Single`.

One may be tempted to make the following definition instead:

```
data S' = Beside' S' S'
      | Single' Char  — too general
```

However, this datatype is too general for the given grammar. An argument of type *Char* can be instantiated to any single character, but we know that this character always has to be `s`. Since there is no choice anyway, there is no extra value in storing that `s`, and the first datatype *S* serves the purpose of encoding the parse trees of the grammar `SequenceOfS` just fine.

For example, the parse tree that corresponds to the first two derivations of the sequence `sss` is represented by the following value of the datatype *S*:

```
Beside Single (Beside Single Single)
```

The third derivation of the sentence `sss` produces the following parse tree:

```
Beside (Beside Single Single) Single
```

To emphasize that these representations contain sufficient information in order to reproduce the original strings, we can write a function that performs this conversion:

```
sToString :: S → String
sToString (Beside l r) = sToString l ++ sToString r
sToString Single      = "s"
```

Applying the function *sToString* to either *Beside Single (Beside Single Single)* or *Beside (Beside Single Single) Single* yields the string `"sss"`.

A *concrete syntax* of a language describes the appearance of the sentences of a language. So the concrete syntax of the language of nonterminal *S* is given by the grammar `SequenceOfS`. concrete syntax

On the other hand, an *abstract syntax* of a language describes the shapes of parse trees of the language, without the need to refer to concrete terminal symbols. Parse trees are therefore often also called abstract syntax trees. The datatype *S* is an example of an abstract syntax for the language of `SequenceOfS`. The adjective ‘abstract’ indicates that values of the abstract syntax do not need to explicitly contain all information about particular sentences, as long as that information is recoverable, as for example by applying function *sToString*. abstract syntax

A function such as *sToString* is often called a *semantic function*. A semantic function is a function that is defined on an abstract syntax of a language. Semantic functions are used to give semantics (meaning) to values. In this example, the meaning of a more abstract representation is expressed in terms of a concrete representation. semantic function

Using the removing left recursion grammar transformation, the grammar `SequenceOfS` can be transformed into the grammar with the following productions:

2. Context-Free Grammars

$$\begin{aligned} S &\rightarrow \mathbf{s}Z \mid \mathbf{s} \\ Z &\rightarrow SZ \mid S \end{aligned}$$

An abstract syntax of this grammar may be given by

$$\begin{aligned} \mathbf{data} \text{ SA} &= \text{ConsS } Z \mid \text{SingleS} \\ \mathbf{data} \text{ Z} &= \text{ConsZ } \text{SA } Z \mid \text{SingleZ } \text{SA} \end{aligned}$$

For each nonterminal in the original grammar, we have introduced a corresponding datatype. For each production for a particular nonterminal (expanding all the alternatives into separate productions), we have introduced a constructor and invented a name for the constructor. The nonterminals on the right hand side of the production rules appear as arguments of the constructors, but the terminals disappear, because that information can be recovered by knowing which constructors have been used.

If we look at the two different grammars for `SequenceOfS`, and the two abstract syntaxes, we can conclude that the only important information about sequences of `s`'s is how many occurrences of `s` there are. So the ultimate abstract syntax for `SequenceOfS` is

$$\mathbf{data} \text{ SS} = \text{Size } \text{Int}$$

Using the abstract syntax `SS`, the sequence `sss` is represented by the parse tree `Size 3`, and we can still recover the original string from a value `SS` by means of a semantic function:

$$\begin{aligned} \text{ssToString} &:: \text{SS} \rightarrow \text{String} \\ \text{ssToString} (\text{Size } n) &= \text{replicate } n \text{ 's'} \end{aligned}$$

The `SequenceOfS` example shows that one may choose between many different abstract syntaxes for a given grammar. The choice of an abstract syntax over another should therefore be determined by the demands of the application, i.e., by what we ultimately want to compute.

Exercise 2.20. The following Haskell datatype represents a limited form of arithmetic expressions

$$\begin{aligned} \mathbf{data} \text{ Expr} &= \text{Add Expr Expr} \\ &\mid \text{Mul Expr Expr} \\ &\mid \text{Con Int} \end{aligned}$$

Give a grammar for a suitable concrete syntax corresponding to this datatype.

Exercise 2.21. Consider the grammar for palindromes that you have constructed in Exercise 2.7. Give parse trees for the palindromes $pal_1 = \text{"abaaba"}$ and $pal_2 = \text{"baaab"}$. Define a datatype `Pal` corresponding to the grammar and represent the parse trees for pal_1 and pal_2 as values of `Pal`.

Exercise 2.22. Consider your answers to Exercises 2.7 and 2.21 where we have given a grammar for palindromes over the alphabet $\{\mathbf{a}, \mathbf{b}\}$ and a Haskell datatype describing the abstract syntax of such palindromes.

1. Write a semantic function that transforms an abstract representation of a palindrome into a concrete one. Test your function with the palindromes pal_1 and pal_2 from Exercise 2.21.
2. Write a semantic function that counts the number of **a**'s occurring in a palindrome. Test your function with the palindromes pal_1 and pal_2 from Exercise 2.21.

Exercise 2.23. Consider your answer to Exercise 2.8, which describes the concrete syntax for mirror palindromes.

1. Define a datatype Mir that describes the abstract syntax corresponding to your grammar. Give the two abstract mirror palindromes $aMir_1$ and $aMir_2$ that correspond to the concrete mirror palindromes $cMir_1 = \text{"abaaba"}$ and $cMir_2 = \text{"abbbba"}$.
2. Write a semantic function that transforms an abstract representation of a mirror palindrome into a concrete one. Test your function with the abstract mirror palindromes $aMir_1$ and $aMir_2$.
3. Write a function that transforms an abstract representation of a mirror palindrome into the corresponding abstract representation of a palindrome (using the datatype from Exercise 2.21). Test your function with the abstract mirror palindromes $aMir_1$ and $aMir_2$.

Exercise 2.24. Consider your answer to Exercise 2.9, which describes the concrete syntax for parity sequences.

1. Define a datatype $Parity$ describing the abstract syntax corresponding to your grammar. Give the two abstract parity sequences $aEven_1$ and $aEven_2$ that correspond to the concrete parity sequences $cEven_1 = \text{"00101"}$ and $cEven_2 = \text{"01010"}$.
2. Write a semantic function that transforms an abstract representation of a parity sequence into a concrete one. Test your function with the abstract parity sequences $aEven_1$ and $aEven_2$.

Exercise 2.25. Consider your answer to Exercise 2.18, which describes the concrete syntax for bit lists by means of a grammar that is not left-recursive.

1. Define a datatype $BitList$ that describes the abstract syntax corresponding to your grammar. Give the two abstract bit-lists $aBitList_1$ and $aBitList_2$ that correspond to the concrete bit-lists $cBitList_1 = \text{"0,1,0"}$ and $cBitList_2 = \text{"0,0,1"}$.
2. Write a semantic function that transforms an abstract representation of a bit list into a concrete one. Test your function with the abstract bit lists $aBitList_1$ and $aBitList_2$.
3. Write a function that concatenates two abstract representations of a bit lists into a bit list. Test your function with the abstract bit lists $aBitList_1$ and $aBitList_2$.

2.7. Constructions on grammars

This section introduces some constructions on grammars that are useful when specifying larger grammars, for example for programming languages. Furthermore, it gives an example of a larger grammar that is transformed in several steps.

The BNF notation, introduced in Section 2.2.1, was first used in the early sixties when the programming language ALGOL 60 was defined and until now it is the

2. Context-Free Grammars

standard way of defining the syntax of programming languages (see, for instance, the Java Language Grammar). You may object that the Java grammar contains more “syntactical sugar” than the grammars that we considered thus far (and to be honest, this also holds for the ALGOL 60 grammar): one encounters nonterminals with postfixes ‘?’, ‘+’ and ‘*’.

EBNF

This extended BNF notation, *EBNF*, is introduced in order to help abbreviate a number of standard constructions that usually occur quite often in the syntax of a programming language:

- one or zero occurrences of nonterminal P , abbreviated $P?$,
- one or more occurrences of nonterminal P , abbreviated P^+ ,
- and zero or more occurrences of nonterminal P , abbreviated P^* .

We could easily express these constructions by adding additional nonterminals, but that decreases the readability of the grammar. The notation for the EBNF constructions is not entirely standardized. In some texts, you will for instance find the notation $[P]$ instead of $P?$, and $\{P\}$ for P^* . The same notation can be used for languages, grammars, and sequences of terminal and nonterminal symbols instead of just single nonterminals. In this section, we define the meaning of these constructs.

We introduced grammars as an alternative for the description of languages. Designing a grammar for a specific language may not be a trivial task. One approach is to decompose the language and to find grammars for each of its constituent parts.

In Definition 2.3, we have defined a number of operations on languages using operations on sets. We now show that these operations can be expressed in terms of operations on context-free grammars.

Theorem 2.10 (Language operations). *Suppose we have grammars for the languages L and M , say $G_L = (T, N_L, R_L, S_L)$ and $G_M = (T, N_M, R_M, S_M)$. We assume that the nonterminal sets N_L and N_M are disjoint. Then*

- the language $L \cup M$ is generated by the grammar (T, N, R, S) where S is a fresh nonterminal, $N = N_L \cup N_M \cup \{S\}$ and $R = R_L \cup R_M \cup \{S \rightarrow S_L, S \rightarrow S_M\}$;
- the language LM is generated by the grammar (T, N, R, S) where S is a fresh nonterminal, $N = N_L \cup N_M \cup \{S\}$ and $R = R_L \cup R_M \cup \{S \rightarrow S_L S_M\}$;
- the language L^* is generated by the grammar (T, N, R, S) where S is a fresh nonterminal, $N = N_L \cup \{S\}$ and $R = R_L \cup \{S \rightarrow \varepsilon, S \rightarrow S_L S\}$;
- the language L^+ is generated by the grammar (T, N, R, S) where S is a fresh nonterminal, $N = N_L \cup \{S\}$ and $R = R_L \cup \{S \rightarrow S_L, S \rightarrow S_L S\}$.

The theorem above establishes that the set-theoretic operations at the level of languages (i. e., sets of sentences) have a direct counterpart at the level of grammatical descriptions. A straightforward question to ask is now: can we also define languages

as the difference between two languages or as the intersection of two languages, and translate these operations to operations on grammars? Unfortunately, the answer is negative – there are no operations on grammars that correspond to the language intersection and difference operators.

Two of the above constructions are important enough to actually define them as grammar operations. Furthermore, we add a new grammar construction for an “optional grammar”.

Definition 2.11 (Grammar operations). Let $G = (T, N, R, S)$ be a context-free grammar and let S' be a fresh nonterminal. Then

$$\begin{aligned} G^* &= (T, N \cup \{S'\}, R \cup \{S' \rightarrow \varepsilon, S' \rightarrow S S'\}, S') \\ G^+ &= (T, N \cup \{S'\}, R \cup \{S' \rightarrow S, S' \rightarrow S S'\}, S') \\ G^? &= (T, N \cup \{S'\}, R \cup \{S' \rightarrow \varepsilon, S' \rightarrow S \quad \}, S') \end{aligned}$$

The definition of $P^?$, P^+ , and P^* for a sequence of symbols P is very similar to the definitions of the operations on grammars. For example, P^* denotes zero or more concatenations of string P , so Dig^* denotes the language consisting of zero or more digits.

Definition 2.12 (EBNF for sequences). Let P be a sequence of nonterminals and terminals, then

$$\begin{aligned} L(P^*) &= L(Z) \quad \text{with} \quad Z \rightarrow \varepsilon \mid PZ \\ L(P^+) &= L(Z) \quad \text{with} \quad Z \rightarrow P \mid PZ \\ L(P^?) &= L(Z) \quad \text{with} \quad Z \rightarrow \varepsilon \mid P \end{aligned}$$

where Z is a new nonterminal in each definition.

Because the concatenation operator for sequences is associative, the operators \cdot^* and \cdot^+ can also be defined symmetrically:

$$\begin{aligned} L(P^*) &= L(Z) \quad \text{with} \quad Z \rightarrow \varepsilon \mid ZP \\ L(P^+) &= L(Z) \quad \text{with} \quad Z \rightarrow P \mid ZP \end{aligned}$$

Many variations are possible on this theme:

$$L(P^* Q) = L(Z) \quad \text{with} \quad Z \rightarrow Q \mid PZ \tag{2.1}$$

or also

$$L(P Q^*) = L(Z) \quad \text{with} \quad Z \rightarrow P \mid ZQ \tag{2.2}$$

2.7.1. SL: an example

To illustrate EBNF and some of the grammar transformations given in the previous section, we give a larger example. The following grammar generates expressions in a very small programming language, called SL.

$$\begin{aligned}
 Expr &\rightarrow \text{if } Expr \text{ then } Expr \text{ else } Expr \\
 Expr &\rightarrow Expr \text{ where } Decls \\
 Expr &\rightarrow AppExpr \\
 AppExpr &\rightarrow AppExpr \text{ } Atomic \mid Atomic \\
 Atomic &\rightarrow Var \mid Number \mid Bool \mid (Expr) \\
 Decls &\rightarrow Decl \\
 Decls &\rightarrow Decls ; Decls \\
 Decl &\rightarrow Var = Expr
 \end{aligned}$$

where the nonterminals *Var*, *Number*, and *Bool* generate variables, number expressions, and boolean expressions, respectively. Note that the brackets around the *Expr* in the production for *Atomic*, and the semicolon in between the *Decl*s in the second production for *Decl*s are also terminal symbols. The following ‘program’ is a sentence of this language:

if true then funny true else false where funny = 7

It is clear that this is not a very convenient language to write programs in.

The above grammar is ambiguous (why?), and we introduce priorities to resolve some of the ambiguities. Application binds stronger than **if**, and both application and **if** bind stronger than **where**. Using the “introduction of priorities” grammar transformation, we obtain:

$$\begin{aligned}
 Expr &\rightarrow Expr_1 \\
 Expr &\rightarrow Expr_1 \text{ where } Decls \\
 Expr_1 &\rightarrow Expr_2 \\
 Expr_1 &\rightarrow \text{if } Expr_1 \text{ then } Expr_1 \text{ else } Expr_1 \\
 Expr_2 &\rightarrow Atomic \\
 Expr_2 &\rightarrow Expr_2 \text{ } Atomic
 \end{aligned}$$

where *Atomic* and *Decl*s have the same productions as before.

The nonterminal *Expr*₂ is left-recursive. Removing left recursion gives the following productions for *Expr*₂:

$$\begin{aligned}
 Expr_2 &\rightarrow Atomic \mid Atomic \text{ } Expr'_2 \\
 Expr'_2 &\rightarrow Atomic \mid Atomic \text{ } Expr'_2
 \end{aligned}$$

Since the new nonterminal $Expr'_2$ has exactly the same productions as $Expr_2$, these productions can be replaced by

$$Expr_2 \rightarrow Atomic \mid Atomic Expr_2$$

So $Expr_2$ generates a nonempty sequence of atomics. Using the \cdot^+ -notation introduced before, we can replace $Expr_2$ by $Atomic^+$.

Another source of ambiguity are the productions for $Decls$. The nonterminal $Decls$ generates a nonempty list of declarations, and the separator $;$ is assumed to be associative. Hence we can apply the “associative separator” transformation to obtain

$$Decls \rightarrow Decl \mid Decls ; Decl$$

or, according to (2.2),

$$Decls \rightarrow Decl (; Decl)^*$$

The last grammar transformation we apply is “left factoring”. This transformation is applied to the productions for $Expr$, and yields

$$\begin{aligned} Expr &\rightarrow Expr_1 Expr'_1 \\ Expr'_1 &\rightarrow \varepsilon \mid \mathbf{where} Decls \end{aligned}$$

Since nonterminal $Expr'_1$ generates either nothing or a **where** clause, we can replace $Expr'_1$ by an optional **where** clause in the production for $Expr$:

$$Expr \rightarrow Expr_1 (\mathbf{where} Decls)?$$

After all these grammar transformations, we obtain the following grammar.

$$\begin{aligned} Expr &\rightarrow Expr_1 (\mathbf{where} Decls)? \\ Expr_1 &\rightarrow Atomic^+ \\ Expr_1 &\rightarrow \mathbf{if} Expr_1 \mathbf{then} Expr_1 \mathbf{else} Expr_1 \\ Atomic &\rightarrow Var \mid Number \mid Bool \mid (Expr) \\ Decls &\rightarrow Decl (; Decl)^* \end{aligned}$$

Exercise 2.26. Give the EBNF notation for each of the basic languages defined in Section 2.3.1.

Exercise 2.27. Let G be a grammar G . Give the language that is generated by $G?$ (i.e., the $\cdot?$ operation applied to G).

Exercise 2.28. Let

$$\begin{aligned} L_1 &= \{ a^m b^m c^n \mid m, n \in \mathbb{N} \} \\ L_2 &= \{ a^m b^n c^n \mid m, n \in \mathbb{N} \} \end{aligned}$$

1. Give grammars for L_1 and L_2 .
2. Is $L_1 \cap L_2$ context-free, i. e., can you give a context-free grammar for this language?

2.8. Parsing

This section formulates the parsing problem, and discusses some of the future topics of the course.

Definition 2.13 (Parsing problem). Given the grammar G and a string s , the *parsing problem* answers the question whether or not $s \in L(G)$. If $s \in L(G)$, the answer to this question may be either a parse tree or a derivation.

parsing problem

This question may not be easy to answer given an arbitrary grammar. Until now we have only seen simple grammars for which it is relatively easy to determine whether or not a string is a sentence of the grammar. For more complicated grammars this may be more difficult. However, in the first part of this course we will show how – given a grammar with certain reasonable properties – we can easily construct parsers by hand. At the same time we will show how the parsing process can quite often be combined with the algorithm we actually want to perform on the recognized object (the semantic function). The techniques we describe comprise a simple, although surprisingly efficient, introduction into the area of compiler construction.

A compiler for a programming language consists of several parts. Examples of such parts are a scanner, a parser, a type checker, and a code generator. Usually, a parser is preceded by a *scanner* (also called *lexer*), which splits an input sentence into a list of so-called tokens. For example, given the sentence

scanner
lexer

```
if true then funny true else false where funny = 7
```

a scanner might return the following list of tokens:

```
["if", "true", "then", "funny", "true",  
 "else", "false", "where", "funny", "=", "7"]
```

token

So a *token* is a syntactical entity. A scanner usually performs the first step towards an abstract syntax: it throws away layout information such as spacing and newlines. In this course we will concentrate on parsers, but some of the concepts of scanners will sometimes be used.

In the second part of this course we will take a look at more complicated grammars, which do not always conform to the restrictions just referred to. By analysing the grammar we may nevertheless be able to generate parsers as well. Such generated parsers will be in such a form that it will be clear that writing such parsers by hand is far from attractive, and actually impossible for all practical cases.

One of the problems we have not referred to yet in this rather formal chapter is of a more practical nature. Quite often the sentence presented to the parser will not be a sentence of the language since mistakes were made when typing the sentence. This raises another interesting question: *What are the minimal changes that have*

to be made to the sentence in order to convert it into a sentence of the language?
 It goes almost without saying that this is an important question to be answered in practice; one would not be very happy with a compiler which, given an erroneous input, would just reply that the “Input could not be recognised”. One of the most important aspects here is to define metric for deciding about the minimality of a change; humans usually make certain mistakes more often than others. A semicolon can easily be forgotten, but the chance that an `if` symbol is missing is far from likely. This is where grammar engineering starts to play a rôle.

Summary

Starting from a simple example, the language of palindromes, we have introduced the concept of a context-free grammar. Associated concepts, such as derivations and parse trees were introduced.

2.9. Exercises

Exercise 2.29. Do there exist languages L such that $\overline{(L^*)} = (\overline{L})^*$?

Exercise 2.30. Give a language L such that $L = L^*$.

Exercise 2.31. Under which circumstances is $L^+ = L^* - \{\varepsilon\}$?

Exercise 2.32. Let L be a language over alphabet $\{a, b, c\}$ such that $L = L^R$. Does L contain only palindromes?

Exercise 2.33. Consider the grammar with productions

$$\begin{aligned} S &\rightarrow AA \\ A &\rightarrow AAA \\ A &\rightarrow a \\ A &\rightarrow bA \\ A &\rightarrow Ab \end{aligned}$$

1. Which terminal strings can be produced by derivations of four or fewer steps?
2. Give at least two distinct derivations for the string `babbab`.
3. For any $m, n, p \geq 0$, describe a derivation of the string `bmabnabp`.

Exercise 2.34. Consider the grammar with productions

$$\begin{aligned} S &\rightarrow aaB \\ A &\rightarrow bBb \\ A &\rightarrow \varepsilon \\ B &\rightarrow Aa \end{aligned}$$

Show that the string `aabbaabba` cannot be derived from S .

2. Context-Free Grammars

Exercise 2.35. Give a grammar for the language

$$L = \{\omega c \omega^R \mid \omega \in \{a, b\}^*\}$$

This language is known as the *center-marked palindromes* language. Give a derivation of the sentence $abcba$.

Exercise 2.36. Describe the language generated by the grammar:

$$\begin{aligned} S &\rightarrow \varepsilon \\ S &\rightarrow A \\ A &\rightarrow aAb \\ A &\rightarrow ab \end{aligned}$$

Can you find another (preferably simpler) grammar for the same language?

Exercise 2.37. Describe the languages generated by the grammars.

$$\begin{aligned} S &\rightarrow \varepsilon \\ S &\rightarrow A \\ A &\rightarrow Aa \\ A &\rightarrow a \end{aligned}$$

and

$$\begin{aligned} S &\rightarrow \varepsilon \\ S &\rightarrow A \\ A &\rightarrow AaA \\ A &\rightarrow a \end{aligned}$$

Can you find other (preferably simpler) grammars for the same languages?

Exercise 2.38. Show that the languages generated by the grammars G_1 , G_2 en G_3 are the same.

$$\begin{array}{lll} G_1 : & G_2 : & G_3 : \\ S \rightarrow \varepsilon & S \rightarrow \varepsilon & S \rightarrow \varepsilon \\ S \rightarrow aS & S \rightarrow Sa & S \rightarrow a \\ & & S \rightarrow SS \end{array}$$

Exercise 2.39. Consider the following property of grammars:

1. the start symbol is the only nonterminal which may have an empty production (a production of the form $X \rightarrow \varepsilon$),
2. the start symbol does not occur in any alternative.

A grammar having this property is called *non-contracting*. The grammar $A \rightarrow aAb \mid \varepsilon$ does not have this property. Give a non-contracting grammar which describes the same language.

Exercise 2.40. Describe the language L of the grammar

$$A \rightarrow AaA \mid a$$

Give a grammar for L that has no left-recursive productions. Give a grammar for L that has no right-recursive productions.

Exercise 2.41. Describe the language L of the grammar

$$X \rightarrow a \mid Xb$$

Give a grammar for L that has no left-recursive productions. Give a grammar for L that has no left-recursive productions and is non-contracting.

Exercise 2.42. Consider the language L of the grammar

$$\begin{aligned} S &\rightarrow T \mid US \\ T &\rightarrow aSa \mid Ua \\ U &\rightarrow S \mid SUT \end{aligned}$$

Give a grammar for L which uses only productions with two or less symbols on the right hand side. Give a grammar for L which uses only two nonterminals.

Exercise 2.43. Give a grammar for the language of all sequences of 0's and 1's which start with a 1 and contain exactly one 0.

Exercise 2.44. Give a grammar for the language consisting of all nonempty sequences of brackets,

$$\{ (,) \}$$

in which the brackets match. An example sentence of the language is $() (()) ()$. Give a derivation for this sentence.

Exercise 2.45. Give a grammar for the language consisting of all nonempty sequences of two kinds of brackets,

$$\{ (,), [,] \}$$

in which the brackets match. An example sentence in this language is $[()] ()$.

Exercise 2.46. This exercise shows an example (attributed to Noam Chomsky) of an ambiguous English sentence. Consider the following grammar for a part of the English language:

$$\begin{aligned} \textit{Sentence} &\rightarrow \textit{Subject Predicate} . \\ \textit{Subject} &\rightarrow \mathbf{they} \\ \textit{Predicate} &\rightarrow \textit{Verb NounPhrase} \\ \textit{Predicate} &\rightarrow \textit{AuxVerb Verb Noun} \\ \textit{Verb} &\rightarrow \mathbf{are} \\ \textit{Verb} &\rightarrow \mathbf{flying} \\ \textit{AuxVerb} &\rightarrow \mathbf{are} \\ \textit{NounPhrase} &\rightarrow \textit{Adjective Noun} \\ \textit{Adjective} &\rightarrow \mathbf{flying} \\ \textit{Noun} &\rightarrow \mathbf{planes} \end{aligned}$$

Give two different leftmost derivations for the sentence

$\mathbf{they\ are\ flying\ planes.}$

2. Context-Free Grammars

Exercise 2.47. Try to find some ambiguous sentences in your own natural language. Here are some ambiguous Dutch sentences seen in the newspapers:

Vliegen met hartafwijking niet gevaarlijk
Jessye Norman kan niet zingen
Alcohol is voor vrouwen schadelijker dan mannen

Exercise 2.48 (•). Is your grammar for Exercise 2.45 unambiguous? If not, find one which *is* unambiguous.

Exercise 2.49. This exercise deals with a grammar that uses unusual terminal and non-terminal symbols. Assume that \odot , \otimes , and \oplus are nonterminals, and the other symbols are terminals.

$\odot \rightarrow \odot \triangle \otimes$
 $\odot \rightarrow \otimes$
 $\otimes \rightarrow \otimes \diamond \oplus$
 $\otimes \rightarrow \oplus$
 $\oplus \rightarrow \clubsuit$
 $\oplus \rightarrow \spadesuit$

Find a derivation for the sentence $\clubsuit \diamond \clubsuit \triangle \spadesuit$.

Exercise 2.50 (•). Prove, using induction, that the grammar G for palindromes in Section 2.2 does indeed generate the language of palindromes.

Exercise 2.51 (••, no answer provided). Prove that the language generated by the grammar of Exercise 2.33 contains all strings over $\{a, b\}$ where the number of a 's is even and greater than zero.

Exercise 2.52 (no answer provided). Consider the natural numbers in unary notation where only the symbol I is used; thus 4 is represented as $IIII$. Write an algorithm that, given a string w of I 's, determines whether or not w is divisible by 7.

Exercise 2.53 (no answer provided). Consider the natural numbers in reverse binary notation; thus 4 is represented as 001 . Write an algorithm that, given a string w of zeros and ones, determines whether or not w is divisible by 7.

Exercise 2.54 (no answer provided). Let w be a string consisting of a 's and b 's only. Write an algorithm that determines whether or not the number of a 's in w equals the number of b 's in w .

3. Parser combinators

Introduction

This chapter is an informal introduction to writing parsers in a lazy functional language using ‘parser combinators’. Parsers can be written using a small set of basic parsing functions, and a number of functions that combine parsers into more complicated parsers. The functions that combine parsers are called parser combinators. The basic parsing functions do not combine parsers, and are therefore not parser combinators in this sense, but they are usually also called parser combinators.

Parser combinators are used to write parsers that are very similar to the grammar of a language. Thus writing a parser amounts to translating a grammar to a functional program, which is often a simple task.

Parser combinators are built by means of standard functional language constructs like higher-order functions, lists, and datatypes. List comprehensions are used in a few places, but they are not essential, and could easily be rephrased using the *map*, *filter* and *concat* functions. Type classes are only used for overloading the equality and arithmetic operators.

We will start by motivating the definition of the type of parser functions. Using that type, we can build parsers for the language of (possibly ambiguous) grammars. Next, we will introduce some elementary parsers that can be used for parsing the terminal symbols of a language.

In Section 3.3 the first parser combinators are introduced, which can be used for sequentially and alternatively combining parsers, and for calculating so-called semantic functions during the parse. Semantic functions are used to give meaning to syntactic structures. As an example, we construct a parser for strings of matching parentheses in Section 3.3.1. Different semantic values are calculated for the matching parentheses: a tree describing the structure, and an integer indicating the nesting depth.

In Section 3.4 we introduce some new parser combinators. Not only do these make life easier later, but their definitions are also nice examples of using parser combinators. A real application is given in Section 3.5, where a parser for arithmetical expressions is developed. Finally, the expression parser is generalised to expressions with an arbitrary number of precedence levels. This is done without coding the priorities of operators as integers, and we will avoid using indices and ellipses.

3. Parser combinators

It is not always possible to directly construct a parser from a context-free grammar using parser combinators. If the grammar is left-recursive, it has to be transformed into a non left-recursive grammar before we can construct a combinator parser. Another limitation of the parser combinator technique as described in this chapter is that it is not trivial to write parsers for complex grammars that perform reasonably efficient. However, there do exist implementations of parser combinators that perform remarkably well, see [10, 12]. For example, there exist good parsers using parser combinators for the Haskell language.

Most of the techniques introduced in this chapter have been described by Burge [4], Wadler [13] and Hutton [7].

This chapter is a revised version of an article by Fokker [5].

Goals

This chapter introduces the first programs for parsing in these lecture notes. Parsers are composed from simple parsers by means of parser combinators. Hence, important primary goals of this chapter are:

- to understand how to parse, i. e., how to recognise structure in a sequence of symbols, by means of parser combinators;
- to be able to construct a parser when given a grammar;
- to understand the concept of semantic functions.

Two secondary goals of this chapter are:

- to develop the capability to abstract;
- to understand the concept of domain specific language.

Required prior knowledge

To understand this chapter, you should be able to formulate the parsing problem, and you should understand the concept of context-free grammar. Furthermore, you should be familiar with functional programming concepts such as **type**, **class**, and higher-order functions.

3.1. The type of parsers

The goals of this section are:

- develop a type *Parser* that is used to give the type of parsing functions;
- show how to obtain this type by means of several abstraction steps.

The *parsing problem* is (see Section 2.8): Given a grammar G and a string s , determine whether or not $s \in L(G)$. If $s \in L(G)$, the answer to this question may be either a parse tree or a derivation. For example, in Section 2.4 we have seen a grammar for sequences of s 's:

$$S \rightarrow SS \mid s$$

A parse tree of an expression of this language is a value of the datatype S (or a value of several variants of that type, see Section 2.6, depending on what you want to do with the result of the parser), which is defined by

$$\mathbf{data} S = \mathit{Beside} S S \mid \mathit{Single}$$

A parser for expressions could be implemented as a function of the following type:

$$\mathbf{type} Parser = String \rightarrow S \quad \text{— preliminary}$$

For parsing substructures, a parser can call other parsers, or call itself recursively. These calls do not only have to communicate their result, but also the part of the input string that is left unprocessed. For example, when parsing the string sss , a parser will first build a parse tree $\mathit{Beside} \mathit{Single} \mathit{Single}$ for ss , and only then build a complete parse tree

$$\mathit{Beside} (\mathit{Beside} \mathit{Single} \mathit{Single}) \mathit{Single}$$

using the unprocessed part s of the input string. As this cannot be done using a global variable, the unprocessed input string has to be part of the result of the parser. The two results can be paired. A better definition for the type *Parser* is hence:

$$\mathbf{type} Parser = String \rightarrow (S, String) \quad \text{— still preliminary}$$

Any parser of type *Parser* returns an S and a *String*. However, for different grammars we want to return different parse trees: the type of tree that is returned depends on the grammar for which we want to parse sentences. Therefore it is better to abstract from the type S , and to turn the parser type into a polymorphic type. The type *Parser* is parametrised with a type a , which represents the type of parse trees.

$$\mathbf{type} Parser a = String \rightarrow (a, String) \quad \text{— still preliminary}$$

For example, a parser that returns a structure of type *Oak* (whatever that is) now has type *Parser Oak*. A parser that parses sequences of s 's has type *Parser S*.

3. Parser combinators

We might also define a parser that does not return a value of type S , but instead the number of s 's in the input sequence. This parser would have type *Parser Int*. Another instance of a parser is a parse function that recognises a string of digits, and returns the number represented by it as a parse 'tree'. In this case the function is also of type *Parser Int*. Finally, a recogniser that either accepts or rejects sentences of a grammar returns a boolean value, and will have type *Parser Bool*.

Until now, we have assumed that every string can be parsed in exactly one way. In general, this need not be the case: it may be that a single string can be parsed in various ways, or that there is no way to parse a string. For example, the string "sss" has the following two parse trees:

Beside (Beside Single Single) Single
Beside Single (Beside Single Single)

As another refinement of the type *Parser*, instead of returning one parse tree (and its associated rest string), we let a parser return a *list* of trees. Each element of the result consists of a tree, paired with the rest string that was left unprocessed after parsing. The type definition of *Parser* therefore becomes:

type *Parser a = String* → [(*a*, *String*)] — useful, but still suboptimal

If there is just one parse, the result of the parse function is a singleton list. If no parse is possible, the result is an empty list. In case of an ambiguous grammar, the result consists of all possible parses.

list of successes

This method for parsing is called the *list of successes* method, described by Wadler [13]. It can be used in situations where in other languages you would use so-called backtracking techniques. In the Bird and Wadler textbook it is used to solve combinatorial problems like the eight queens problem [3]. If only one solution is required rather than all possible solutions, you can take the head of the list of successes. Thanks to lazy evaluation, not all elements of the list are computed if only the first value is needed, so there will be no loss of efficiency. Lazy evaluation provides a backtracking approach to finding the first solution.

lazy evaluation

Parsers with the type described so far operate on strings, that is lists of characters. There is however no reason for not allowing parsing strings of elements other than characters. You may imagine a situation in which a preprocessor prepares a list of tokens (see Section 2.8), which is subsequently parsed. To cater for this situation we refine the parser type once more: we let the type of the elements of the input string be an argument of the parser type. Calling the type of symbols s , and as before the result type a , the type of parsers becomes

type *Parser s a = [s]* → [(*a*, [*s*])]

or, if you prefer meaningful identifiers over conciseness:

— The type of parsers

```
type Parser symbol result = [symbol] → [(result, [symbol])]
```

Listing 3.1: ParserType.hs

```
type Parser symbol result = [symbol] → [(result, [symbol])]
```

We will use this type definition in the rest of this chapter. This type is defined in Listing 3.1, the first part of our parser library. The list of successes appears in result type of a parser. Each element of this list is a possible parsing of (an initial part of) the input. We will hardly use the full generality provided by the *Parser* type: the type of the input *s* (or *symbol*) will almost always be *Char*.

3.2. Elementary parsers

The goals of this section are:

- introduce some very simple parsers for parsing sentences of grammars with rules of the form:

$$\begin{aligned} A &\rightarrow \varepsilon \\ A &\rightarrow \mathbf{a} \\ A &\rightarrow x \end{aligned}$$

where *x* is a sequence of terminals;

- show how one can construct useful functions from simple, trivially correct functions by means of generalisation and partial parametrisation.

This section defines parsers that can only be used to parse fixed sequences of terminal symbols. For a grammar with a production that contains nonterminals in its right-hand side we need techniques that will be introduced in the following section.

We will start with a very simple parse function that just recognises the terminal symbol **a**. The type of the input string symbols is *Char* in this case, and as a parse ‘tree’ we also simply use a *Char*:

```
symbola :: Parser Char Char
symbola [] = []
symbola (x : xs) | x == 'a' = [('a', xs)]
                  | otherwise = []
```

3. Parser combinators

— Elementary parsers

```
symbol :: Eq s => s -> Parser s s
symbol a [] = []
symbol a (x : xs) | x == a = [(x, xs)]
                  | otherwise = []

satisfy :: (s -> Bool) -> Parser s s
satisfy p [] = []
satisfy p (x : xs) | p x = [(x, xs)]
                  | otherwise = []

token :: Eq s => [s] -> Parser s [s]
token k xs | k == take n xs = [(k, drop n xs)]
          | otherwise = []
    where n = length k

failp :: Parser s a
failp xs = []

succeed :: a -> Parser s a
succeed r xs = [(r, xs)]
```

— Applications of elementary parsers

```
digit :: Parser Char Char
digit = satisfy isDigit
```

Listing 3.2: ParserType.hs

The list of successes method immediately pays off, because now we can return an empty list if no parsing is possible (because the input is empty, or does not start with an `a`).

In the same fashion, we can write parsers that recognise other symbols. As always, rather than defining a lot of closely related functions, it is better to abstract from the symbol to be recognised by making it an extra argument of the function. Furthermore, the function can operate on lists of characters, but also on lists of symbols of other types, so that it can be used in other applications than character oriented ones. The only prerequisite is that the symbols to be parsed can be tested for equality. In Haskell, this is indicated by the *Eq* predicate in the type of the function.

Using these generalisations, we obtain the function *symbol* that is given in Listing 3.2. The function *symbol* is a function that, given a symbol, returns a parser for that symbol. A parser on its turn is a function too. This is why two arguments appear in the definition of *symbol*.

We will now define some elementary parsers that can do the work traditionally taken

care of by lexical analysers (see Section 2.8). For example, a useful parser is one that recognises a fixed string of symbols, such as `while` or `switch`. We will call this function *token*; it is defined in Listing 3.2. As in the case of the *symbol* function we have parametrised this function with the string to be recognised, effectively making it into a family of functions. Of course, this function is not confined to strings of characters. However, we do need an equality test on the type of values in the input string; the type of *token* is:

$$\textit{token} :: \textit{Eq} \ s \Rightarrow [s] \rightarrow \textit{Parser} \ s \ [s]$$

The function *token* is a generalisation of the *symbol* function, in that it recognises a list of symbols instead of a single symbol. Note that we cannot define *symbol* in terms of *token*: the two functions have incompatible types.

Another generalisation of *symbol* is a function which may, depending on the input, return different parse results. Instead of specifying a specific symbol, we can parametrise the function with a *condition* that the symbol should fulfill. Thus the function *satisfy* has a function $s \rightarrow \textit{Bool}$ as argument. Where *symbol* tests for equality to a specific value, the function *satisfy* tests for compliance with this predicate. It is defined in Listing 3.2. This generalised function is for example useful when we want to parse digits (characters in between '0' and '9'):

$$\begin{aligned} \textit{digit} &:: \textit{Parser} \ \textit{Char} \ \textit{Char} \\ \textit{digit} &= \textit{satisfy} \ \textit{isDigit} \end{aligned}$$

where the function *isDigit* is the standard predicate that tests whether or not a character is a digit:

$$\begin{aligned} \textit{isDigit} &:: \textit{Char} \rightarrow \textit{Bool} \\ \textit{isDigit} \ x &= '0' \leq x \wedge x \leq '9' \end{aligned}$$

In books on grammar theory an empty string is often called ' ϵ '. In this tradition, we will define a function *epsilon* that 'parses' the empty string. It does not consume any input, and hence always returns an empty parse tree and unmodified input. A zero-tuple can be used as a result value: $()$ is the only value of the type $()$ – both the type and the value are pronounced *unit*.

unit type

$$\begin{aligned} \textit{epsilon} &:: \textit{Parser} \ s \ () \\ \textit{epsilon} \ xs &= [(), xs] \end{aligned}$$

A more useful variant is the function *succeed*, which also doesn't consume input, but always returns a given, fixed value (or 'parse tree', if you can call the result of processing zero symbols a parse tree). It is defined in Listing 3.2.

Dual to the function *succeed* is the function *failp*, which fails to recognise any symbol on the input string. As the result list of a parser is a 'list of successes', and in the case of failure there are no successes, the result list should be empty. Therefore the

3. Parser combinators

function *failp* always returns the empty list of successes. It is defined in Listing 3.2. Note the difference with *epsilon*, which *does* have one element in its list of successes (albeit an empty one).

Do not confuse *failp* with *epsilon*: there is an important difference between returning one solution (which contains the unchanged input as ‘rest’ string) and not returning a solution at all!

Exercise 3.1. Define a function *capital* :: *Parser Char Char* that parses capital letters.

Exercise 3.2. Since *satisfy* is a generalisation of *symbol*, the function *symbol* can be defined as an instance of *satisfy*. How can this be done?

Exercise 3.3. Define the function *epsilon* using *succeed*.

3.3. Parser combinators

Using the elementary parsers from the previous section, parsers can be constructed for terminal symbols from a grammar. More interesting are parsers for *nonterminal* symbols. It is convenient to *construct* these parsers by partially parametrising higher-order functions.

The goals of this section are:

- show how parsers can be constructed directly from the productions of a grammar. The kind of productions for which parsers will be constructed are

$$\begin{aligned} A &\rightarrow x \mid y \\ A &\rightarrow x y \end{aligned}$$

where *x* and *y* are sequences of nonterminal or terminal symbols;

- show how we can construct a small, powerful combinator language (a domain specific language) for the purpose of parsing;
- understand and use the concept of semantic functions in parsers.

Let us have a look at the grammar for expressions again, see also Section 2.5:

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

where *Digs* is a nonterminal that generates the language of sequences of digits, see Section 2.3.1. An expression can be parsed according to any of the two rules for *E*. This implies that we want to have a way to say that a parser consists of several alternative parsers. Furthermore, the first rule says that in order to parse an expression, we should first parse a term, then a terminal symbol *+*, and then an expression. This

implies that we want to have a way to say that a parser consists of several parsers that are applied sequentially.

So important operations on parsers are sequential and alternative composition: a more complex construct can consist of a simple construct *followed* by another construct (sequential composition), or by a *choice* between two constructs (alternative composition). These operations correspond directly to their grammatical counterparts. We will develop two functions for this, which for notational convenience are defined as operators: $\langle * \rangle$ for sequential composition, and $\langle | \rangle$ for alternative composition. The names of these operators are chosen so that they can be easily remembered: $\langle * \rangle$ ‘multiplies’ two constructs together, and $\langle | \rangle$ can be pronounced as ‘or’. Be careful, though, not to confuse the $\langle | \rangle$ -operator with Haskell’s built-in construct `|`, which is used to distinguish cases in a function definition or as a separator for the constructors of a datatype.

Priorities of these operators are defined so as to minimise parentheses in practical situations:

```
infixl 6  $\langle * \rangle$ 
infixr 4  $\langle | \rangle$ 
```

So $\langle * \rangle$ has a higher priority – i. e., it binds stronger – than $\langle | \rangle$.

Both operators take two parsers as argument, and return a parser as result. By again combining the result with other parsers, you may construct even more involved parsers.

In the definitions in Listing 3.3, the functions operate on parsers p and q . Apart from the arguments p and q , the function operates on a string, which can be thought of as the string that is parsed by the parser that is the result of combining p and q .

We start with the definition of operator $\langle * \rangle$. For sequential composition, p must be applied to the input first. After that, q is applied to the rest string of the result. The first parser, p , returns a list of successes, each of which contains a value and a rest string. The second parser, q , should be applied to the rest string, returning a second value. Therefore we use a list comprehension, in which the second parser is applied in all possible ways to the rest string of the first parser:

$$(p \langle * \rangle q) \text{ } xs = [(combine \ r_1 \ r_2, \ zs) \\ | (r_1, \ ys) \leftarrow p \ \text{ } xs \\ , (r_2, \ zs) \leftarrow q \ \text{ } ys \\]$$

The rest string of the parser for the sequential composition of p and q is whatever the second parser q leaves behind as rest string.

3. Parser combinators

— Parser combinators

```
(<|>)      :: Parser s a → Parser s a → Parser s a
(p <|> q) xs = p xs ++ q xs

(<*>)      :: Parser s (b → a) → Parser s b → Parser s a
(p <*> q) xs = [(f x, zs)
                | (f , ys) ← p xs
                  , ( x, zs) ← q ys
                ]

(<$>)      :: (a → b) → Parser s a → Parser s b
(f <$> p) xs = [(f y, ys)
                | ( y, ys) ← p xs
                ]
```

— Applications of parser combinators

```
newdigit :: Parser Char Int
newdigit = f <$> digit
  where f c = ord c - ord '0'
```

Listing 3.3: ParserCombinators.hs

Now, how should the results of the two parsings be combined? We could, of course, parametrise the whole thing with an operator that describes how to combine the parts (as is done in the *zipWith* function). However, we have chosen for a different approach, which nicely exploits the ability of functional languages to manipulate functions. The function *combine* should combine the results of the two parse trees recognised by *p* and *q*. In the past, we have interpreted the word ‘tree’ liberally: simple values, like characters, may also be used as a parse ‘tree’. We will now also accept *functions* as parse trees. That is, the result type of a parser may be a function type.

If the first parser that is combined by *<*>* would return a function of type $b \rightarrow a$, and the second parser a value of type *b*, a straightforward choice for the *combine* function would be function application. That is exactly the approach taken in the definition of *<*>* in Listing 3.3. The first parser returns a function, the second parser a value, and the combined parser returns the value that is obtained by applying the function to the value.

Apart from ‘sequential composition’ we need a parser combinator for representing ‘choice’. For this, we have the parser combinator operator *<|>*. Thanks to the list of successes method, both *p*₁ and *p*₂ return lists of possible parsings. To obtain all possible parsings when applying *p*₁ or *p*₂, we only need to concatenate these two lists.

By combining parsers with parser combinators we can construct new parsers. The most important parser combinators are `<*>` and `<|>`. The parser combinator `<,>` in exercise 3.12 is just a variation of `<*>`.

Sometimes we are not quite satisfied with the result value of a parser. The parser might work well in that it consumes symbols from the input adequately (leaving the unused symbols as rest-string in the tuples in the list of successes), but the result value might need some postprocessing. For example, a parser that recognises one digit is defined using the function *satisfy*: `digit = satisfy isDigit`. In some applications, we may need a parser that recognises one digit character, but returns the result as an integer, instead of a character. In a case like this, we can use a new parser combinator: `<$>`. It takes a function and a parser as argument; the result is a parser that recognises the same string as the original parser, but ‘postprocesses’ the result using the function. We use the \$ sign in the name of the combinator, because the combinator resembles the operator that is used for normal function application in Haskell: `f $ x = f x`. The definition of `<$>` is given in Listing 3.3. It is an infix operator:

```
infixl 7 <$>
```

Using this postprocessing parser combinator, we can modify the parser *digit* that was defined above:

```
newdigit :: Parser Char Int
newdigit = f <$> digit
  where f c = ord c - ord '0'
```

The auxiliary function *f* determines the ordinal number of a digit character; using the parser combinator `<$>` it is applied to the result part of the *digit* parser.

In practice, the `<$>` operator is used to build a certain value during parsing (in the case of parsing a computer program this value may be the generated code, or a list of all variables with their types, etc.). Put more generally: using `<$>` we can add *semantic functions* to parsers.

A parser for the `SequenceOfS` grammar that returns the abstract syntax tree of the input, i.e., a value of type *S*, see Section 2.6, is defined as follows:

```
sequenceOfS :: Parser Char S
sequenceOfS = Beside <$> sequenceOfS <*> sequenceOfS
  <|> const Single <$> symbol 's'
```

But if you try to run this function, you will get a stack overflow! If you apply *sequenceOfS* to a string, the first thing it does is to apply itself to the same string, which loops. The problem stems from the fact that the underlying grammar is left-recursive. For any left-recursive grammar, a systematically constructed parser using parser combinators will exhibit the problem that it loops. However, in Section 2.5

3. Parser combinators

we have shown how to remove the left recursion in the `SequenceOfS` grammar. The resulting grammar is used to obtain the following parser:

```
sequenceOfS' :: Parser Char SA
sequenceOfS' =
    const ConsS <$> symbol 's' <*> parseZ
  <|> const SingleS <$> symbol 's'
  where parseZ = ConsZ <$> sequenceOfS' <*> parseZ
        <|> SingleZ <$> sequenceOfS'
```

This example is a direct translation of the grammar obtained by using the removing left recursion grammar transformation. There exists a much simpler parser for parsing sequences of `s`'s.

Exercise 3.4. Prove for all $f :: a \rightarrow b$ that

$$f \text{ < \$ > } \text{ succeed } a = \text{ succeed } (f \ a)$$

In the sequel we will often use this rule for constant functions f , i. e., $f = \lambda_ \rightarrow c$ for some term c .

Exercise 3.5. Consider the parser $(:) \text{ < \$ > } \text{ symbol 'a'}$. Give its type and show its results on inputs `[]` and $x : xs$.

Exercise 3.6. Consider the parser $(:) \text{ < \$ > } \text{ symbol 'a' < * > } p$. Give its type and show its results on inputs `[]` and $x : xs$.

Exercise 3.7. Define a parser for Booleans.

Exercise 3.8 (no answer provided). Define parsers for each of the basic languages defined in Section 2.3.1.

Exercise 3.9. Consider the grammar for palindromes that you have constructed in Exercise 2.7.

1. Give the datatype Pal_2 that corresponds to this grammar.
2. Define a parser $palin_2$ that returns parse trees for palindromes. Test your function with the palindromes $cPal_1 = \text{"abaaba"}$ and $cPal_2 = \text{"baaab"}$. Compare the results with your answer to Exercise 2.21.
3. Define a parser $palina$ that counts the number of `a`'s occurring in a palindrome.

Exercise 3.10. Consider the grammar for a part of the English language that is given in Exercise 2.46.

1. Give the datatype $English$ that corresponds to this grammar.
2. Define a parser $english$ that returns parse trees for the English language. Test your function with the sentence `they are flying planes`. Compare the result to your answer of Exercise 2.46.

Exercise 3.11. When defining the priority of the $\langle| \rangle$ operator with the **infix** keyword, we also specified that the operator associates to the right. Why is this a better choice than association to the left?

Exercise 3.12. Define a parser combinator \langle, \rangle that combines two parsers. The value returned by the combined parser is a tuple containing the results of the two component parsers. What is the type of this parser combinator?

Exercise 3.13. The term ‘parser combinator’ is in fact not an adequate description for $\langle \$ \rangle$. Can you think of a better word?

Exercise 3.14. Compare the type of $\langle \$ \rangle$ with the type of the standard function *map*. Can you describe your observations in an easy-to-remember, catchy phrase?

Exercise 3.15. Define $\langle * \rangle$ in terms of \langle, \rangle and $\langle \$ \rangle$. Define \langle, \rangle in terms of $\langle * \rangle$ and $\langle \$ \rangle$.

Exercise 3.16. If you examine the definitions of $\langle * \rangle$ and $\langle \$ \rangle$ in Listing 3.3, you can observe that $\langle \$ \rangle$ is in a sense a special case of $\langle * \rangle$. Can you define $\langle \$ \rangle$ in terms of $\langle * \rangle$?

3.3.1. Matching parentheses: an example

Using parser combinators, it is often fairly straightforward to construct a parser for a language for which you have a grammar. Consider, for example, the grammar that you wrote in Exercise 2.44:

$$S \rightarrow (S) S \mid \varepsilon$$

This grammar can be directly translated to a parser, using the parser combinators $\langle * \rangle$ and $\langle | \rangle$. We use $\langle * \rangle$ when symbols are written next to each other, and $\langle | \rangle$ when $|$ appears in a production (or when there is more than one production for a nonterminal).

```
parens :: Parser Char ??? — ill-typed
parens = symbol '(' <*> parens <*> symbol ')' <*> parens
        <|> epsilon
```

However, this function is not correctly typed: the parsers in the first alternative cannot be composed using $\langle * \rangle$, as for example *symbol '('* is not a parser returning a function.

But we can postprocess the parser *symbol '('* so that, instead of a character, this parser *does* return a function. So, what function should we use? This depends on the kind of value that we want as a result of the parser. A nice result would be a tree-like description of the parentheses that are parsed. For this purpose we introduce an abstract syntax, see Section 2.6, for the parentheses grammar. We obtain the following Haskell datatype:

```
data Parentheses = Match Parentheses Parentheses
                 | Empty
```

3. Parser combinators

```
data Parentheses = Match Parentheses Parentheses
                  | Empty
deriving Show
open = symbol '('
close = symbol ')'
parens :: Parser Char Parentheses
parens =          (λ_ x _ y → Match x y)
               <$> open <*> parens <*> close <*> parens
               <|> succeed Empty
nesting :: Parser Char Int
nesting =         (λ_ x _ y → max (1 + x) y)
               <$> open <*> nesting <*> close <*> nesting
               <|> succeed 0
```

Listing 3.4: ParseParentheses.hs

For example, the sentence `()()` is represented by

```
Match Empty (Match Empty Empty)
```

Suppose we want to calculate the number of parentheses in a sentence. The number of parentheses is calculated by the function `nrofpars`, which is defined by induction on the datatype `Parentheses`.

```
nrofpars :: Parentheses → Int
nrofpars (Match pl pr) = 2 + nrofpars pl + nrofpars pr
nrofpars Empty         = 0
```

Using the datatype `Parentheses`, we can add ‘semantic functions’ to the parser. We then obtain the definition of `parens` in Listing 3.4.

By varying the function used as a first argument of `<$>` (the ‘semantic function’), we can return other things than parse trees. As an example we construct a parser that calculates the nesting depth of nested parentheses, see the function `nesting` defined in Listing 3.4.

A session in which `nesting` is used may look like this:

```
? nesting "()()()"
[(2,[]), (2,"()"), (1,"()()"), (0,"()()()")]
? nesting "()"
[(1,""), (0,"()")]
```


As you can see, when there is a syntax error in the argument, there are no solutions with empty rest string. It is fairly simple to test whether a given string belongs to the language that is parsed by a given parser.

Exercise 3.17. What is the type of the function $f = \lambda x y \rightarrow Match\ x\ y$ which appears in function *parens* in Listing 3.4? What is the type of the parser *open*? Using the type of $\langle \$ \rangle$, what is the type of $f\ \langle \$ \rangle\ open$? Can $f\ \langle \$ \rangle\ open$ be used as a left hand side of $\langle * \rangle parens$? What is the type of the result?

Exercise 3.18. What is a convenient way for $\langle * \rangle$ to associate? Does it?

Exercise 3.19. Write a function *test* that determines whether or not a given string belongs to the language parsed by a given parser.

3.4. More parser combinators

In principle you can build parsers for any context-free language using the combinators $\langle * \rangle$ and $\langle | \rangle$, but in practice it is easier to have some more parser combinators available. In traditional grammar formalisms, additional symbols are used to describe for example optional or repeated constructions. Consider for example the BNF formalism, in which originally only sequential and alternative composition can be used (denoted by juxtaposition and vertical bars, respectively), but which was later extended to EBNF to also allow for repetition, denoted by a star. The goal of this section is to show how the set of parser combinators can be extended.

3.4.1. Parser combinators for EBNF

It is very easy to make new parser combinators for EBNF. As a first example we consider repetition. Given a parser p for a construction, *many* p constructs a parser for zero or more occurrences of that construction:

$$\begin{aligned} \text{many} &:: \text{Parser } s\ a \rightarrow \text{Parser } s\ [a] \\ \text{many } p &= (\cdot) \langle \$ \rangle p \langle * \rangle \text{many } p \\ &\langle | \rangle \text{ succeed } [] \end{aligned}$$

So the EBNF expression P^* is implemented by *many* P . The function (\cdot) is just the cons-operator for lists: it takes a head element and a tail list and combines them.

The order in which the alternatives are given only influences the order in which solutions are placed in the list of successes.

For example, the *many* combinator can be used in parsing a natural number:

$$\begin{aligned} \text{natural} &:: \text{Parser } \text{Char } \text{Int} \\ \text{natural} &= \text{foldl } f\ 0 \langle \$ \rangle \text{many } \text{newdigit} \\ &\text{where } f\ a\ b = a * 10 + b \end{aligned}$$

3. Parser combinators

— EBNF parser combinators

option :: *Parser s a* → *a* → *Parser s a*

option p d = *p* <|> *succeed d*

many :: *Parser s a* → *Parser s [a]*

many p = (*:*) <\$> *p* <*> *many p* <|> *succeed []*

many₁ :: *Parser s a* → *Parser s [a]*

many₁ p = (*:*) <\$> *p* <*> *many p*

pack :: *Parser s a* → *Parser s b* → *Parser s c* → *Parser s b*

pack p r q = ($\lambda_ x _ \rightarrow x$) <\$> *p* <*> *r* <*> *q*

listOf :: *Parser s a* → *Parser s b* → *Parser s [a]*

listOf p s = (*:*) <\$> *p* <*> *many ((_ x → x) <\$> s <*> p)*

— Auxiliary functions

first :: *Parser s b* → *Parser s b*

first p xs | *null r* = []
| *otherwise* = [*head r*]

where *r* = *p xs*

greedy, greedy₁ :: *Parser s b* → *Parser s [b]*

greedy = *first . many*

greedy₁ = *first . many₁*

Listing 3.5: EBNF.hs

Defined in this way, the *natural* parser also accepts empty input as a number. If this is not desired, we had better use the *many*₁ parser combinator, which accepts one or more occurrences of a construction, and corresponds to the EBNF expression P^+ , see Section 2.7. It is defined in Listing 3.5. Another combinator from EBNF is the *option* combinator $P?$. It takes a parser as argument, and returns a parser that recognises the same construct, but which also succeeds if that construct is not present in the input string. The definition is given in Listing 3.5. It has an additional argument: the value that should be used as result in case the construct is not present. It is a kind of ‘default’ value.

By the use of the *option* and *many* functions, a large amount of backtracking possibilities are introduced. This is not always advantageous. For example, if we define a parser for identifiers by

$$identifier = many_1 (satisfy isAlpha)$$

a single word may also be parsed as two identifiers. Caused by the order of the alternatives in the definition of *many* (*succeed* [] appears as the second alternative), the ‘greedy’ parsing, which accumulates as many letters as possible in the identifier is tried first, but if parsing fails elsewhere in the sentence, also less greedy parsings of the identifier are tried – in vain. You will give a better definition of *identifier* in Exercise 3.27.

In situations where from the way the grammar is built we can predict that it is hopeless to try non-greedy results of *many*, we can define a parser transformer *first*, that transforms a parser into a parser that only returns the first possible parsing. It does so by taking the first element of the list of successes.

$$\begin{aligned} first &:: Parser\ a\ b \rightarrow Parser\ a\ b \\ first\ p\ xs\ | \ null\ r &= [] \\ &| \ otherwise = [head\ r] \\ \mathbf{where}\ r &= p\ xs \end{aligned}$$

Using this function, we can create a special ‘take all or nothing’ version of *many*:

$$\begin{aligned} greedy &= first . many \\ greedy_1 &= first . many_1 \end{aligned}$$

If we compose the *first* function with the *option* parser combinator:

$$obligatory\ p\ d = first\ (option\ p\ d)$$

we get a parser which must accept a construction if it is present, but which does not fail if it is not present.

3. Parser combinators

3.4.2. Separators

The combinators *many*, *many₁* and *option* are classical in compiler constructions – there are notations for it in EBNF (\cdot^* , \cdot^+ and $\cdot?$, respectively) –, but there is no need to leave it at that. For example, in many languages constructions are frequently enclosed between two meaningless symbols, most often some sort of parentheses. For this case we design a parser combinator *pack*. Given a parser for an opening token, a body, and a closing token, it constructs a parser for the enclosed body, as defined in Listing 3.5. Special cases of this combinator are:

```
parenthesised p = pack (symbol '(') p (symbol ')')
bracketed p     = pack (symbol '[') p (symbol ']')
compound p     = pack (token "begin") p (token "end")
```

Another frequently occurring construction is repetition of a certain construction, where the elements are separated by some symbol. You may think of lists of arguments (expressions separated by commas), or compound statements (statements separated by semicolons). For the parse trees, the separators are of no importance. The function *listOf* below generates a parser for a non-empty list, given a parser for the items and a parser for the separators:

```
listOf :: Parser s a → Parser s b → Parser s [a]
listOf p s = (:) <$> p <*> many ((λ_ x → x) <$> s <*> p)
```

Useful instantiations are:

```
commaList, semicList :: Parser Char a → Parser Char [a]
commaList p = listOf p (symbol ',' )
semicList p = listOf p (symbol ';' )
```

A somewhat more complicated variant of the function *listOf* is the case where the separators carry a meaning themselves. For example, in arithmetical expressions, where the operators that separate the subexpressions have to be part of the parse tree. For this case we will develop the functions *chainr* and *chainl*. These functions expect that the parser for the separators returns a function (!); that function is used by *chain* to combine parse trees for the items. In the case of *chainr* the operator is applied right-to-left, in the case of *chainl* it is applied left-to-right. The functions *chainr* and *chainl* are defined in Listing 3.6 (remember that $\$$ is function application: $f \$ x = f x$).

The definitions look quite complicated, but when you look at the underlying grammar they are quite straightforward. Suppose we apply operator \oplus (\oplus is an operator variable, it denotes an arbitrary right-associative operator) from right to left, so

$$e_1 \oplus e_2 \oplus e_3 \oplus e_4$$

=

— Chain expression combinators

```

chainr :: Parser s a → Parser s (a → a → a) → Parser s a
chainr pe po = h <$> many (j <$> pe <*> po) <*> pe
  where j x op = (x `op` )
        h fs x = foldr ($) x fs

chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl pe po = h <$> pe <*> many (j <$> po <*> pe)
  where j op x = (op `x` )
        h x fs = foldl (flip ($)) x fs

```

Listing 3.6: Chains.hs

$$\begin{aligned}
& e_1 \oplus (e_2 \oplus (e_3 \oplus e_4)) \\
= & \\
& ((e_1 \oplus) \cdot (e_2 \oplus) \cdot (e_3 \oplus)) e_4
\end{aligned}$$

It follows that we can parse such expressions by parsing many pairs of expressions and operators, turning them into functions, and applying all those functions to the last expression. This is done by function *chainr*, see Listing 3.6.

If operator \oplus is applied from left to right, then

$$\begin{aligned}
& e_1 \oplus e_2 \oplus e_3 \oplus e_4 \\
= & \\
& ((e_1 \oplus e_2) \oplus e_3) \oplus e_4 \\
= & \\
& ((\oplus e_4) \cdot (\oplus e_3) \cdot (\oplus e_2)) e_1
\end{aligned}$$

So such an expression can be parsed by first parsing a single expression (e_1), and then parsing many pairs of operators and expressions, turning them into functions, and applying all those functions to the first expression. This is done by function *chainl*, see Listing 3.6.

Functions *chainl* and *chainr* can be made more efficient by avoiding the construction of the intermediate list of functions. The resulting definitions can be found in the article by Fokker [5].

Note that functions *chainl* and *chainr* are very similar, the only difference is that everything is ‘turned around’: function *j* of *chainr* takes a value and an operator, and returns the function obtained by ‘left’ applying the operator; function *j* of *chainl* takes an operator and a value, and returns the function obtained by ‘right’ applying the operator to the value. Such functions are sometimes called *dual*.

dual

3. Parser combinators

Exercise 3.20.

1. Define a parser that analyses a string and recognises a list of digits separated by a space character. The result is a list of integers.
2. Define a parser *sumParser* that recognises digits separated by the character '+' and returns the sum of these integers.
3. Both parsers return a list of solutions. What should be changed in order to get only one solution?

Exercise 3.21. What is the value of

$$\text{many (symbol 'a') xs}$$

for $xs \in \{\ [], ['a'], ['b'], ['a', 'b'], ['a', 'a', 'b']\}$?

Exercise 3.22. Consider the application of the parser *many (symbol 'a')* to the string `aaa`. In what order do the four possible parsings appear in the list of successes?

Exercise 3.23 (no answer provided). Using the parser combinators *option*, *many* and *many₁* define parsers for each of the basic languages defined in Section 2.3.1.

Exercise 3.24. As another variation on the theme 'repetition', define a parser combinator *psequence* that transforms a *list of parsers* for some type into a *parser returning a list* of elements of that type. What is the type of *psequence*? Also define a combinator *choice* that iterates the operator `<|>`.

Exercise 3.25. As an application of *psequence*, define the function *token* that was discussed in Section 3.2.

Exercise 3.26 (no answer provided). Carefully analyse the semantic functions in the definition of *chainl* in Listing 3.6.

Exercise 3.27. In real programming languages, identifiers follow rather flexible rules: the first symbol must be a letter, but the symbols that follow (if any) may be a letter, digit, or underscore symbol. Define a more realistic parser *identifier*.

3.5. Arithmetical expressions

The goal of this section is to use parser combinators in a concrete application. We will develop a parser for arithmetical expressions, which have the following concrete syntax:

$$\begin{array}{l} E \quad \rightarrow E + E \\ \quad \quad | E - E \\ \quad \quad | E / E \\ \quad \quad | (E) \\ \quad \quad | Digs \end{array}$$

Besides these productions, we also have productions for identifiers and applications of functions:

$$\begin{aligned}
 E &\rightarrow \textit{Identifier} \\
 &\quad | \textit{Identifier} (\textit{Args}) \\
 \textit{Args} &\rightarrow \varepsilon \mid E (, E)^*
 \end{aligned}$$

The parse trees for this grammar are of type *Expr*:

```

data Expr = Con Int
           | Var String
           | Fun String [Expr]
           | Expr :+: Expr
           | Expr :-: Expr
           | Expr **: Expr
           | Expr :/: Expr

```

You can almost recognise the structure of the parser in this type definition. But in order to account for the priorities of the operators, we will use a grammar with three non-terminals ‘expression’, ‘term’ and ‘factor’: an expression is composed of terms separated by + or −; a term is composed of factors separated by * or /, and a factor is a constant, variable, function call, or expression between parentheses.

This grammar appears as a parser in the functions in Listing 3.7.

The first parser, *fact*, parses factors.

```

fact :: Parser Char Expr
fact = Con <$> integer
      <|> Var <$> identifier
      <|> Fun <$> identifier <*> parenthesised (commaList expr)
      <|> parenthesised expr

```

The first alternative is an integer parser which is postprocessed by the ‘semantic function’ *Con*. The second and third alternative are a variable or function call, depending on the presence of an argument list. In absence of the latter, the function *Var* is applied, in presence the function *Fun*. For the fourth alternative there is no semantic function, because the meaning of an expression between parentheses is the meaning of the expression.

For the definition of a term as a list of factors separated by multiplicative operators we use the function *chainl*. Recall that *chainl* repeatedly recognises its first argument (*fact*), separated by its second argument (a * or /). The parse trees for the individual factors are joined by the constructor functions that appear before <\$>. We use *chainl* and not *chainr* because the operator ‘/’ is considered to be left-associative.

The function *expr* is analogous to *term*, only with additive operators instead of multiplicative operators, and with *terms* instead of *factors*.

```

— Type definition for parse tree
data Expr = Con Int
          | Var String
          | Fun String [Expr]
          | Expr :+: Expr
          | Expr :-: Expr
          | Expr :*: Expr
          | Expr :/: Expr

— Parser for expressions with two priorities
fact :: Parser Char Expr
fact =   Con <$> integer
        <|> Var <$> identifier
        <|> Fun <$> identifier <*> parenthesised (commaList expr)
        <|> parenthesised expr

integer :: Parser Char Int
integer = (const negate <$> (symbol '-')) 'option' id <*> natural

term :: Parser Char Expr
term = chainl fact
      (
        const (:*) <$> symbol '*'
        <|> const (:/) <$> symbol '/'
      )

expr :: Parser Char Expr
expr = chainl term
      (
        const (:+) <$> symbol '+'
        <|> const (-:) <$> symbol '-'
      )

```

Listing 3.7: ExpressionParser.hs

This example clearly shows the strength of parsing with parser combinators. There is no need for a separate formalism for grammars; the production rules of the grammar are combined with higher-order functions. Also, there is no need for a separate parser generator (like ‘yacc’); the functions can be viewed both as description of the grammar and as an executable parser.

Exercise 3.28. 1. Give the parse tree for the expressions "abc", "(abc)", "a*b+1", "a*(b+1)", "-1-a", and "a(1,b)".

2. Why is the parse tree for the expression "a(1,b)" not the first solution of the parser? Modify the functions in Listing 3.7 in such way that it will be.

Exercise 3.29. A function with no arguments such as "f()" is not accepted by the parser. Explain why and modify the parser in such way that it will be.

Exercise 3.30. Modify the functions in Listing 3.7, in such a way that + is parsed as a right-associative operator, and - is parsed as a left-associative operator.

3.6. Generalised expressions

This section generalises the parser in the previous section with respect to priorities. Arithmetical expressions in which operators have more than two levels of priority can be parsed by writing more auxiliary functions between *term* and *expr*. The function *chainl* is used in each definition, with as first argument the function of one priority level lower.

If there are nine levels of priority, we obtain nine copies of almost the same text. This is not as it should be. Functions that resemble each other are an indication that we should write a generalised function, where the differences are described using extra arguments. Therefore, let us inspect the differences in the definitions of *term* and *expr* again. These are:

- The operators and associated tree constructors that are used in the second argument of *chainl*
- The parser that is used as first argument of *chainl*

The generalised function will take these two differences as extra arguments: the first in the form of a list of pairs, the second in the form of a parse function:

```

type Op a = (Char, a → a → a)
gen :: [Op a] → Parser Char a → Parser Char a
gen ops p = chainl p (choice (map f ops))
  where f (s, c) = const c <$> symbol s

```

If furthermore we define as shorthand:

3. Parser combinators

```
multis = [( '*', (:*)), (' / ', (: / ))]
addis  = [( '+', (:+)), (' - ', (:-))]
```

then *expr* and *term* can be defined as partial parametrisations of *gen*:

```
expr = gen addis term
term = gen multis fact
```

By expanding the definition of *term* in that of *expr* we obtain:

```
expr = addis 'gen' (multis 'gen' fact)
```

which an experienced functional programmer immediately recognises as an application of *foldr*:

```
expr = foldr gen fact [addis, multis]
```

From this definition a generalisation to more levels of priority is simply a matter of extending the list of operator-lists.

The very compact formulation of the parser for expressions with an arbitrary number of priority levels is possible because the parser combinators can be used together with the existing mechanisms for generalisation and partial parametrisation in Haskell.

Contrary to conventional approaches, the levels of priority need not be coded explicitly with integers. The only thing that matters is the relative position of an operator in the list of 'list with operators of the same priority'. Also, the insertion of new priority levels is very easy. The definitions are summarised in Listing 3.8.

Summary

This chapter shows how to construct parsers from simple combinators. It shows how a small parser combinator library can be a powerful tool in the construction of parsers. Furthermore, this chapter gives a rather basic implementation of the parser combinator library. More advanced implementations are discussed elsewhere.

3.7. Exercises

Exercise 3.31. How should the parser of Section 3.6 be adapted to also allow raising an expression to the power of an expression?

Exercise 3.32. Prove the following laws

$$h \langle \$ \rangle (f \langle \$ \rangle p) = (h . f) \langle \$ \rangle p \tag{3.1}$$

$$h \langle \$ \rangle (p \langle | \rangle q) = (h \langle \$ \rangle p) \langle | \rangle (h \langle \$ \rangle q) \tag{3.2}$$

$$h \langle \$ \rangle (p \langle * \rangle q) = ((h.) \langle \$ \rangle p) \langle * \rangle q \tag{3.3}$$

— Parser for expressions with arbitrary many priorities

```

type Op a = (Char, a → a → a)
fact' :: Parser Char Expr
fact' = Con <$> integer
      <|> Var <$> identifier
      <|> Fun <$> identifier <*> parenthesised (commaList expr')
      <|> parenthesised expr'

gen :: [Op a] → Parser Char a → Parser Char a
gen ops p = chainl p (choice (map f ops))
  where f (s, c) = const c <$> symbol s

expr' :: Parser Char Expr
expr' = foldr gen fact' [addis, multis]
multis = [( ' * ' , (:*:) ), ( ' / ' , (:/:) )]
addis  = [( ' + ' , (:+:) ), ( ' - ' , (:-) )]

```

Listing 3.8: GExpressionParser.hs

Exercise 3.33. Consider your answer to Exercise 2.23. Define a combinator parser $pMir$ that transforms a concrete representation of a mirror-palindrome into an abstract one. Test your function with the concrete mirror-palindromes $cMir_1$ and $cMir_2$.

Exercise 3.34. Consider your answer to Exercise 2.25. Assuming the comma is an associative operator, we can give the following abstract syntax for bit-lists:

```

data BitList = SingleB Bit | ConsB Bit BitList

```

Define a combinator parser $pBitList$ that transforms a concrete representation of a bit-list into an abstract one. Test your function with the concrete bit-lists $cBitList_1$ and $cBitList_2$.

Exercise 3.35. Define a parser for fixed-point numbers, that is numbers like 12.34 and -123.456. Also integers are acceptable. Notice that the part following the decimal point looks like an integer, but has a different semantics!

Exercise 3.36. Define a parser for floating point numbers, which are fixed point numbers followed by an optional E and an (positive or negative, integer) exponent.

Exercise 3.37. Define a parser for Java assignments that consist of a variable, an = sign, an expression and a semicolon.

Exercise 3.38 (no answer provided). Define a parser for (simplified) Java statements.

Exercise 3.39 (no answer provided). Outline the construction of a parser for Java programs.

3. Parser combinators

4. Grammar and Parser design

The previous chapters have introduced many concepts related to grammars and parsers. The goal of this chapter is to review these concepts, and to show how they are used in the design of grammars and parsers.

The design of a grammar and parser for a language consists of several steps: you have to

1. give example sentences of the language for which you want to design a grammar and a parser;
2. give a grammar for the language for which you want to have a parser;
3. test that the grammar can indeed describe the example sentences;
4. analyse this grammar to find out whether or not it has some desirable properties;
5. possibly transform the grammar to obtain some of these desirable properties;
6. decide on the type of the parser: *Parser a b*, that is, decide on both the input type *a* of the parser (which may be the result type of a scanner), and the result type *b* of the parser.
7. construct a basic parser;
8. add semantic functions;
9. test that the parser can parse the example sentences you have given in the first step, and that the parser returns what you expect.

We will describe and exemplify each of these steps in detail in the rest of this section.

As a running example we will construct a grammar and parser for travelling schemes for day trips, of the following form:

Groningen 8:37 9:44 Zwolle 9:49 10:15 Utrecht 10:21 11:05 Den Haag

We might want to do several things with such a schema, for example:

1. compute the net travel time, i.e., the travel time minus the waiting time (2 hours and 17 minutes in the above example);

4. Grammar and Parser design

2. compute the total time one has to wait on the intermediate stations (11 minutes).

This chapter defines functions to perform these computations.

4.1. Step 1: Example sentences for the language

We have already given an example sentence above:

Groningen 8:37 9:44 Zwolle 9:49 10:15 Utrecht 10:21 11:05 Den Haag

Other example sentences are:

Utrecht Centraal 10:25 10:58 Amsterdam Centraal
Assen

4.2. Step 2: A grammar for the language

The starting point for designing a parser for your language is to define a grammar that describes the language as precisely as possible. It is important to convince yourself from the fact that the grammar you give really generates the desired language, since the grammar will be the basis for grammar transformations, which might turn the grammar into a set of incomprehensible productions.

For the language of travelling schemes, we can give several grammars. The following grammar focuses on the fact that a trip consists of zero or more departures and arrivals.

$$\begin{aligned} TS &\rightarrow TS \textit{ Departure Arrival } TS \mid \textit{ Station} \\ \textit{ Station} &\rightarrow \textit{ Identifier}^+ \\ \textit{ Departure} &\rightarrow \textit{ Time} \\ \textit{ Arrival} &\rightarrow \textit{ Time} \\ \textit{ Time} &\rightarrow \textit{ Nat} : \textit{ Nat} \end{aligned}$$

where *Identifier* and *Nat* have been defined in Section 2.3.1. So a travelling scheme is a sequence of departure and arrival times, separated by stations. Note that a single station is also a travelling scheme with this grammar.

Another grammar focuses on changing at a station:

$$\begin{aligned} TS &\rightarrow \textit{ Station Departure (Arrival Station Departure)^* Arrival Station} \\ &\mid \textit{ Station} \end{aligned}$$

So each travelling scheme starts and ends at a station, and in between there is a list of intermediate stations.

4.3. Step 3: Testing the grammar

Both grammars we have given in step 2 describe the example sentences given in step 1. The derivation of these sentences using these grammars is easy.

4.4. Step 4: Analysing the grammar

To parse sentences of a language efficiently, we want to have a unambiguous grammar that is left-factored and not left recursive. Depending on the parser we want to obtain, we might desire other properties of our grammar. So a first step in designing a parser is analysing the grammar, and determining which properties are (not) satisfied. We have not yet developed tools for grammar analysis (we will do so in the chapter on *LL(1)* parsing) but for some grammars it is easy to detect some properties.

The first example grammar is left and right recursive: the first production for *TS* starts and ends with *TS*. Furthermore, the sequence *Departure Arrival* is an associative separator in the generated language.

These properties may be used for transforming the grammar. Since we don't mind about right recursion, we will not make use of the fact that the grammar is right recursive. The other properties will be used in grammar transformations in the following subsection.

4.5. Step 5: Transforming the grammar

Since the sequence *Departure Arrival* is an associative separator in the generated language, the productions for *TS* may be transformed into:

$$TS \rightarrow Station \mid Station\ Departure\ Arrival\ TS \quad (4.1)$$

Thus we have removed the left recursion in the grammar. Both productions for *TS* start with the nonterminal *Station*, so *TS* can be left factored. The resulting productions are:

$$\begin{aligned} TS &\rightarrow Station\ Z \\ Z &\rightarrow \varepsilon \mid Departure\ Arrival\ TS \end{aligned}$$

We can also apply equivalence (2.1) to the two productions for *TS* from (4.1), and obtain the following single production:

$$TS \rightarrow (Station\ Departure\ Arrival)^* Station \quad (4.2)$$

So which productions do we take for *TS*? This depends on what we want to do with the parsed sentences. We will show several choices in the next section.

4.6. Step 6: Deciding on the types

We want to write a parser for travel schemes, that is, we want to write a function ts of type

$$ts :: \text{Parser } ??$$

The question marks should be replaced by the input type and the result type, respectively. For the input type we can choose between at least two possibilities: characters, Char or tokens Token . The type of tokens can be chosen as follows:

```
data Token = Station-Token Station | Time-Token Time
type Station = String
type Time = (Int, Int)
```

We will construct a parser for both input types in the next subsection. So ts has one of the following two types.

$$ts :: \text{Parser } \text{Char } ?$$

$$ts :: \text{Parser } \text{Token } ?$$

For the result type we have many choices. If we just want to compute the total travelling time, Int suffices for the result type. If we want to compute the total travelling time, the total waiting time, and a nicely printed version of the travelling scheme, we may do several things:

- define three parsers, with Int (total travelling time), Int (total waiting time), and String (nicely printed version) as result type, respectively;
- define a single parser with the triple $(\text{Int}, \text{Int}, \text{String})$ as result type;
- define an abstract syntax for travelling schemes, say a datatype TS , and define three functions on TS that compute the desired results.

The first alternative parses the input three times, and is rather inefficient compared with the other alternatives. The second alternative is hard to extend if we want to compute something extra, but in some cases it might be more efficient than the third alternative. The third alternative needs an abstract syntax. There are several ways to define an abstract syntax for travelling schemes. The first abstract syntax corresponds to definition (4.1) of grammar TS .

```
data TS1 = Single1 Station
          | Cons1 Station Time Time TS1
```

where Station and Time are defined above. A second abstract syntax corresponds to the grammar for travelling schemes defined in (4.2).

```
type TS2 = [(Station, Time, Time), Station]
```


So a travelling scheme is a tuple, the first component of which is a list of triples consisting of a departure station, a departure time, and an arrival time, and the second component of which is the final arrival station. A third abstract syntax corresponds to the second grammar defined in Section 4.2:

```
data  $TS_3 = Single_3 Station$ 
      |  $Cons_3 (Station, Time, [(Time, Station, Time)], Time, Station)$ 
```

Which abstract syntax should we take? Again, this depends on what we want to do with the abstract syntax. Since TS_2 and TS_1 combine departure and arrival times in a tuple, they are convenient to use when computing travelling times. TS_3 is useful when we want to compute waiting times since it combines arrival and departure times in one constructor. Often we want to exactly mimic the productions of the grammar in the abstract syntax, so if we use grammar (4.1) for travelling schemes, we use TS_1 for the abstract syntax. Note that TS_1 is a datatype, whereas TS_2 is a type. TS_1 cannot be defined as a type because of the two alternative productions for TS . TS_2 can be defined as a datatype by adding a constructor. Types and datatypes each have their advantages and disadvantages; the application determines which to use. The result type of the parsing function ts may be one of types mentioned earlier (Int , etc.), or one of TS_1 , TS_2 , TS_3 .

4.7. Step 7: Constructing the basic parser

Converting a grammar to a parser is a mechanical process that consists of a set of simple replacement rules. Functional programming languages offer some extra flexibility that we sometimes use, but usually writing a parser is a simple translation. We use the following replacement rules.

grammar construct	Haskell/parser construct
\rightarrow	$=$
$ $	$\langle \rangle$
(space)	$\langle * \rangle$
$\cdot +$	$many_1$
$\cdot *$	$many$
$\cdot ?$	$option$
terminal x	$symbol\ x$
begin of sequence of symbols	$undefined\langle \$ \rangle$

Note that we start each sequence of symbols by $undefined\langle \$ \rangle$. The $undefined$ has to be replaced by an appropriate semantic function in Step 6, but putting $undefined$ here ensures type correctness of the parser. Of course, running the parser will result in an error.

4. Grammar and Parser design

We construct a basic parser for each of the input types *Char* and *Token*.

4.7.1. Basic parsers from strings

Applying these rules to the grammar (4.2) for travelling schemes, we obtain the following basic parser.

```
station :: Parser Char Station
station = undefined <$> many1 identifier

time :: Parser Char Time
time = undefined <$> natural <*> symbol ':' <*> natural

departure, arrival :: Parser Char Time
departure = undefined <$> time
arrival   = undefined <$> time

tsstring :: Parser Char ?
tsstring = undefined
          <$> many (
                    undefined
                    <$> spaces
                    <*> station
                    <*> spaces
                    <*> departure
                    <*> spaces
                    <*> arrival
                  )
          <*> spaces
          <*> station

spaces :: Parser Char String
spaces = undefined <$> many (symbol ' ')
```

The only thing left to do is to add the semantic glue to the functions. The semantic glue also determines the type of the function *tsstring*, which is denoted by ? for the moment. For the other basic parsers we have chosen some reasonable return types. The semantic functions are defined in the next and final step.

4.7.2. A basic parser from tokens

To obtain a basic parser from tokens, we first write a scanner that produces a list of tokens from a string.

```
scanner :: String → [Token]
scanner = mkTokens . combine . words

combine :: [String] → [String]
```

```

combine []           = []
combine [x]         = [x]
combine (x : y : xs) = if isAlpha (head x) ^ isAlpha (head y)
                        then combine ((x ++ " " ++ y) : xs)
                        else x : combine (y : xs)

mkToken :: String → Token
mkToken xs = if isDigit (head xs)
              then Time-Token (mkTime xs)
              else Station-Token xs

parse_result :: [(a, b)] → a
parse_result xs
  | null xs    = error "parse_result: could not parse the input"
  | otherwise  = fst (head xs)

mkTime :: String → Time
mkTime = parse_result . time

```

This is a basic scanner with very basic error messages, but it suffices for now. The composition of the scanner with the function *tstoken₁* defined below gives the final parser.

```

tstoken1 :: Parser Token ?
tstoken1 = undefined
          <$> many ( undefined
                    <$> tstation
                    <*> tdeparture
                    <*> tarrival
                    )
          <*> tstation

tstation :: Parser Token Station
tstation (Station-Token s : xs) = [(s, xs)]
tstation _                      = []

tdeparture, tarrival :: Parser Token Time
tdeparture (Time-Token (h, m) : xs) = [(h, m), xs]
tdeparture _                        = []
tarrival (Time-Token (h, m) : xs) = [(h, m), xs]
tarrival _                          = []

```

where again the semantic functions remain to be defined. Note that functions *tdeparture* and *tarrival* are the same functions. Their presence reflects their presence in the grammar.

Another basic parser from tokens is based on the second grammar of Section 4.2.

```

tstoken2 :: Parser Token ?
tstoken2 = undefined

```

4. Grammar and Parser design

```
<$> tstation
<*> tdeparture
<*> many ( undefined
           <$> tarrival
           <*> tstation
           <*> tdeparture
         )
<*> tarrival
<*> tstation
<|> undefined <$> tstation
```

4.8. Step 8: Adding semantic functions

Once we have the basic parsing functions, we need to add the semantic glue: the functions that take the results of the elements in the right hand side of a production, and convert them into the result of the left hand side. The basic rule is: Let the types do the work!

First we add semantic functions to the basic parsing functions *station*, *time*, *departure*, *arrival*, and *spaces*. Since function *many₁ identifier* returns a list of strings, and we want to obtain the concatenation of these strings for the station name, we can take the concatenation function *concat* for *undefined* in function *station*. To obtain a value of type *Time* from an integer, a character, and an integer, we have to combine the two integers in a tuple. So we take the following function

$$\lambda x _ y \rightarrow (x, y)$$

for *undefined* in *time*. Now, since function *time* returns a value of type *Time*, we can take the identity function for *undefined* in *departure* and *arrival*, and then we replace *id <\$> time* by just *time*. Finally, the result of *many* is a string, so for *undefined* in *spaces* we can take the identity function too.

The first semantic function for the basic parser *tsstring* defined in Section 4.7.1 returns an abstract syntax tree of type TS_2 . So the first *undefined* in *tsstring* should return a tuple of a list of things of the correct type (the first component of the type TS_2) and a *Station*. Since *many* returns a list of things, we can construct such a tuple by means of the function

$$\lambda x _ y \rightarrow (x, y)$$

provided *many* returns a value of the desired type: $[(Station, Time, Time)]$. Note that this semantic function basically only throws away the value returned by the *spaces* parser: we are not interested in the spaces between the components of our

travelling scheme. The *many* parser returns a value of the correct type if we replace the second occurrence of *undefined* in *tsstring* by the function

$$\lambda_ x _ y _ z \rightarrow (x, y, z)$$

Again, the results of *spaces* are thrown away. This completes a parser for travelling schemes. The next semantic functions we define compute the net travel time. To compute the net travel time, we have to compute the travel time of each trip from a station to a station, and to add the travel times of all of these trips. We obtain the travel time of a single trip if we replace the second occurrence of *undefined* in *tsstring* by:

$$\lambda_ _ _ (xh, xm) _ (zh, zm) \rightarrow (zh - xh) * 60 + zm - xm$$

and Haskell's prelude function *sum* sums these times, so for the first occurrence of *undefined* we take:

$$\lambda x _ _ \rightarrow sum\ x$$

The final set of semantic functions we define are used for computing the total waiting time. Since the second grammar of Section 4.2 combines arrival times and departure times, we use a parser based on this grammar: the basic parser *tstoken₂*. We have to give definitions of the three *undefined* semantic functions. If a trip consists of a single station, there is now waiting time, so the last occurrence of *undefined* is the function *const 0*. The second occurrence of function *undefined* computes the waiting time for one intermediate station:

$$\lambda(uh, um) _ (wh, wm) \rightarrow (wh - uh) * 60 + wm - um$$

Finally, the first occurrence of *undefined* sums the list of waiting time obtained by means of the function that replaces the second occurrence of *undefined*:

$$\lambda_ _ x _ _ \rightarrow sum\ x$$

4.9. Step 9: Did you get what you expected

In the last step you test your parser(s) to see whether or not you have obtained what you expected, and whether or not you have made errors in the above process.

Summary

This chapter describes the different steps that have to be considered in the design of a grammar and a language.

4. Grammar and Parser design

4.10. Exercises

Exercise 4.1. Write a parser *floatLiteral* for Java float-literals. The EBNF grammar for float-literals is given by:

$$\begin{aligned} \textit{FloatLiteral} &\rightarrow \textit{IntPart} . \textit{FractPart}? \textit{ExponentPart}? \textit{FloatSuffix}? \\ &\quad | . \textit{FractPart} \textit{ExponentPart}? \textit{FloatSuffix}? \\ &\quad | \textit{IntPart} \textit{ExponentPart} \textit{FloatSuffix}? \\ &\quad | \textit{IntPart} \textit{ExponentPart}? \textit{FloatSuffix} \\ \textit{IntPart} &\rightarrow \textit{SignedInteger} \\ \textit{FractPart} &\rightarrow \textit{Digits} \\ \textit{ExponentPart} &\rightarrow \textit{ExponentIndicator} \textit{SignedInteger} \\ \textit{SignedInteger} &\rightarrow \textit{Sign}? \textit{Digits} \\ \textit{Digits} &\rightarrow \textit{Digits} \textit{Digit} \mid \textit{Digit} \\ \textit{ExponentIndicator} &\rightarrow \textit{e} \mid \textit{E} \\ \textit{Sign} &\rightarrow + \mid - \\ \textit{FloatSuffix} &\rightarrow \textit{f} \mid \textit{F} \mid \textit{d} \mid \textit{D} \end{aligned}$$

To keep your parser simple, assume that all nonterminals, except for the nonterminal *FloatLiteral*, are represented by a *String* in the abstract syntax.

Exercise 4.2. Write an evaluator *signedFloat* for Java float-literals (the float-suffix may be ignored).

Exercise 4.3. Up to the definition of the semantic functions, parsers constructed on a (fixed) abstract syntax have the same shape. Give this parsing scheme for Java float literals.

5. Compositionality

Introduction

Many recursive functions follow a common pattern of recursion. These common patterns of recursion can conveniently be captured by higher order functions. For example: many recursive functions defined on lists are instances of the higher order function *foldr*. It is possible to define a function such as *foldr* for a whole range of datatypes other than lists. Such functions are called compositional. Compositional functions on datatypes are defined in terms of algebras of semantic actions that correspond to the constructors of the datatype. Compositional functions can typically be used to define the semantics of programming languages constructs. Such semantics is referred to as algebraic semantics. Algebraic semantics often uses algebras that contain functions from tuples to tuples. Such functions can be seen as computations that read values from a component of the domain and write values to a component of the codomain. The former values are called inherited attributes and the latter values are called synthesised attributes. Attributes can be both inherited and synthesised. As explained in Section 2.6, there is an important relationship between grammars and compositionality: with every grammar, which describes the concrete syntax of a language, one can associate a (possibly mutually recursive) datatype, which describes the abstract syntax of the language. Compositional functions on these datatypes are called syntax driven.

Goals

After studying this chapter and making the exercises you will

- know how to generalise constructors of a datatype to an algebra;
- know how to write compositional functions, also known as folds, on (possibly mutually recursive) datatypes;
- understand the advantages of using folds, and have seen that many problems can be solved with a fold;
- know that a fold applied to the constructors algebra is the identity function;
- have seen the notions of fusion and deforestation;
- know how to write syntax driven code;
- understand the connection between datatypes, abstract syntax and concrete syntax;
- understand the notions of synthesised and inherited attributes;

5. Compositionality

- be able to associate inherited and synthesised attributes with the different alternatives of (possibly mutually recursive) datatypes (or the different nonterminals of grammars);
- be able to define algebraic semantics in terms of compositional (or syntax driven) code that is defined using algebras of computations which make use of inherited and synthesised attributes.

Organisation

The chapter is organised as follows. Section 5.1 shows how to define compositional recursive functions on built-in lists using a function which is similar to *foldr* and shows how to do the same thing with user-defined lists and streams. Section 5.2 shows how to define compositional recursive functions on several kinds of trees. Section 5.3 defines algebraic semantics. Section 5.4 shows the usefulness of algebraic semantics by presenting an expression evaluator, an expression interpreter which makes use of a stack and expression compiler to a stack machine. They only differ in the way they handle basic expressions (variables and local definitions are handled in the same way). All three examples use an algebra of computations which can read values from and write values to an environment which binds names to values. In a second version of the expression evaluator the use of inherited and synthesised attributes is made more explicit by using tuples. Section 5.5 presents a relatively complex example of the use of tuples in combination with compositionality. It deals with the problem of variable scope when compiling block structured languages.

5.1. Lists

This section introduces compositional functions on the well known datatype of lists. Compositional functions are defined on the built-in datatype `[a]` for lists (Section 5.1.1), on a user-defined datatype `List a` for lists (Section 5.1.2), and on streams or infinite lists (Section 5.1.3). We also show how to construct an algebra that directly corresponds to a datatype.

5.1.1. Built-in lists

The datatype of lists is perhaps the most important example of a datatype. A list is either the empty list `[]` or a nonempty list `(x : xs)` consisting of an element `x` at the head of the list, followed by a tail `xs` which itself is again a list. Thus, the type `[x]` is recursive. In Haskell the type `[x]` for lists is built-in. The informal definition of above corresponds to the following (pseudo) datatype.

```
data [x] = x : [x] | []
```


Many recursive functions on lists look very similar. Computing the sum of the elements of a list of integers (*sumL*, where the *L* denotes that it is a sum function defined on lists) is very similar to computing the product of the elements of a list of integers (*prodL*).

$$\begin{aligned} \text{sumL, prodL} &:: [\text{Int}] \rightarrow \text{Int} \\ \text{sumL } (x : xs) &= x + \text{sumL } xs \\ \text{sumL } [] &= 0 \\ \text{prodL } (x : xs) &= x * \text{prodL } xs \\ \text{prodL } [] &= 1 \end{aligned}$$

The function *sumL* replaces the list constructor (*:*) by (+) and the list constructor [] by 0 and the function *prodL* replaces the list constructor (*:*) by (*) and the list constructor [] by 1. Note that we have replaced the constructor (*:*) (a constructor with two arguments) by binary operators (+) and (*) (i.e. functions with two arguments) and the constructor [] (a constructor with zero variables) by constants 0 and 1 (i.e. ‘functions’ with zero variables). The similarity between the definitions of the functions *sumL* and *prodL* can be captured by the following higher order recursive function *foldL*, which is nothing else but an uncurried version of the well known function *foldr*. Don’t confuse *foldL* with Haskell’s prelude function *foldl*, which works the other way around.

$$\begin{aligned} \text{foldL} &:: (x \rightarrow l \rightarrow l, l) \rightarrow [x] \rightarrow l \\ \text{foldL } (op, c) &= \text{fold} \\ &\mathbf{where} \\ &\text{fold } (x : xs) = op \ x \ (\text{fold } xs) \\ &\text{fold } [] = c \end{aligned}$$

The function *foldL* recursively replaces the constructor (*:*) by an operator *op* and the constructor [] by a constant *c*. We can now use *foldL* to compute the sum and product of the elements of a list of integers as follows.

```
? foldL ((+),0) [1,2,3,4]
10
? foldL ((*),1) [1,2,3,4]
24
```

The pair (*op, c*) is often referred to as a list-algebra. More precisely, a list-algebra consists of a type *l* (the carrier of the algebra), a binary operator *op* of type *x* → *l* → *l* and a constant *c* of type *l*. Note that a type (like *Int*) can be the carrier of a list-algebra in more than one way (for example using ((+),0) and ((*),1)). Here is another example of how to turn *Int* into a list-algebra.

```
? foldL (\_ n -> n+1,0) [1,2,3,4]
4
```

5. Compositionality

This list-algebra ignores the value at the head of a list, and increments the result obtained thus far with one. It corresponds to the function *sizeL* defined by:

```
sizeL :: [x] → Int
sizeL (_ : xs) = 1 + sizeL xs
sizeL []       = 0
```

Note that the type of *sizeL* is more general than the types of *sumL* and *prodL*. The type of the elements of the list does not play a role.

5.1.2. User-defined lists

In this subsection we present an example of a fold function defined on another datatype than built-in lists. To keep things simple we redo the list example for user-defined lists.

```
data List x = Cons x (List x) | Nil
```

User-defined lists are defined in the same way as built-in ones. The constructors *(:)* and *[]* are replaced by constructors *Cons* and *Nil*. Here are the types of the constructors *Cons* and *Nil*.

```
? :t Cons
Cons :: a -> List a -> List a
? :t Nil
Nil  :: List a
```

A algebra type *ListAlgebra* corresponding to the datatype *List* directly follows the structure of that datatype.

```
type ListAlgebra x l = (x → l → l, l)
```

The left hand side of the type definition is obtained from the left hand side of the datatype as follows: a postfix *Algebra* is added at the end of the name *List* and a type variable *l* is added at the end of the whole left hand side of the type definition. The right hand side of the type definition is obtained from the right hand side of the data definition as follows: all *List x* valued constructors are replaced by *l* valued functions which have the same number of arguments (if any) as the corresponding constructors. The types of recursive constructor arguments (i.e. arguments of type *List x*) are replaced by recursive function arguments (i.e. arguments of type *l*). The types of the other arguments are simply left unchanged. In a similar way, the definition of a fold function can be generated automatically from the data definition.

```
foldList :: ListAlgebra x l → List x → l
foldList (cons, nil) = fold
```

where

$$\begin{aligned} \text{fold } (\text{Cons } x \text{ } xs) &= \text{cons } x \text{ } (\text{fold } xs) \\ \text{fold } \text{Nil} &= \text{nil} \end{aligned}$$

The constructors *Cons* and *Nil* in the left hand sides of the definition of the local function *fold* are replaced by functions *cons* and *nil* in the right hand sides. The function *fold* is applied recursively to all recursive constructor arguments of type *List x* to return a value of type *l* as required by the functions of the algebra (in this case *Cons* and *cons* have one such recursive argument). The other arguments are left unchanged. Recursive functions on user-defined lists which are defined by means of *foldList* are called compositional. Every algebra defines a unique compositional function. Here are three examples of compositional functions. They correspond to the examples of Section 5.1.1.

$$\begin{aligned} \text{sumList}, \text{prodList} &:: \text{List Int} \rightarrow \text{Int} \\ \text{sumList} &= \text{foldList } ((+), 0) \\ \text{prodList} &= \text{foldList } ((*), 1) \\ \text{sizeList} &:: \text{List } x \rightarrow \text{Int} \\ \text{sizeList} &= \text{foldList } (\text{const } (1+), 0) \end{aligned}$$

It is worth mentioning one particular *ListAlgebra*: the trivial *ListAlgebra* that replaces *Cons* by *Cons* and *Nil* by *Nil*. This algebra defines the identity function on user-defined lists.

$$\begin{aligned} \text{idListAlgebra} &:: \text{ListAlgebra } x \text{ } (\text{List } x) \\ \text{idListAlgebra} &= (\text{Cons}, \text{Nil}) \\ \text{idList} &:: \text{List } x \rightarrow \text{List } x \\ \text{idList} &= \text{foldList } \text{idListAlgebra} \end{aligned}$$

```
? idList (Cons 1 (Cons 2 Nil))
Cons 1 (Cons 2 Nil)
```

5.1.3. Streams

In this section we consider streams (or infinite lists).

streams

```
data Stream x = And x (Stream x)
```

Here is a standard example of a stream: the infinite list of fibonacci numbers.

$$\begin{aligned} \text{fibStream} &:: \text{Stream Int} \\ \text{fibStream} &= \text{And } 0 \text{ } (\text{And } 1 \text{ } (\text{restOf } \text{fibStream})) \\ \text{where} \\ \text{restOf } (\text{And } x \text{ } \text{stream}@(\text{And } y \text{ } -)) &= \text{And } (x + y) \text{ } (\text{restOf } \text{stream}) \end{aligned}$$

5. Compositionality

The algebra type *StreamAlgebra*, and the fold function *foldStream* can be generated automatically from the datatype *Stream*.

```
type StreamAlgebra x s = x → s → s
foldStream :: StreamAlgebra x s → Stream x → s
foldStream and = fold
where
  fold (And x xs) = and x (fold xs)
```

Note that the algebra has only one component because *Stream* has only one constructor. For the same reason the fold function is defined using only one equation. Here is an example of using a compositional function on user defined streams. It computes the first element of a monotone stream that is greater or equal than a given value.

```
firstGreaterThan :: Ord x ⇒ x → Stream x → x
firstGreaterThan n = foldStream (\x y → if x ≥ n then x else y)
```

5.2. Trees

Now that we have seen how to generalise the *foldr* function on built-in lists to compositional functions on user-defined lists and streams we proceed by explaining another common class of datatypes: trees. We will treat four different kinds of trees in the subsections below:

- binary trees;
- trees for matching parentheses;
- expression trees;
- general trees.

Furthermore, we will briefly mention the concepts of fusion and deforestation.

5.2.1. Binary trees

A binary tree is either a node where the tree splits into two subtrees or a leaf which holds a value.

```
data BinTree x = Bin (BinTree x) (BinTree x) | Leaf x
```

One can generate the corresponding algebra type *BinTreeAlgebra* and fold function *foldBinTree* from the datatype automatically. Note that *Bin* has two recursive arguments and that *Leaf* has one non-recursive argument.

```
type BinTreeAlgebra x t = (t → t → t, x → t)
```

```

foldBinTree :: BinTreeAlgebra x t → BinTree x → t
foldBinTree (bin, leaf) = fold
  where
    fold (Bin l r) = bin (fold l) (fold r)
    fold (Leaf x) = leaf x

```

In the *BinTreeAlgebra* type, the *bin* part of the algebra has two arguments of type *t* and the *leaf* part of the algebra has one argument of type *x*. Similarly, in the *foldBinTree* function, the local *fold* function is applied recursively to both arguments of *bin* and is not called on the argument of *leaf*. We can now define compositional functions on binary trees much in the same way as we defined them on lists. Here is an example: the function *sizeBinTree* computes the size of a binary tree.

```

sizeBinTree :: BinTree x → Int
sizeBinTree = foldBinTree ((+), const 1)

? sizeBinTree (Bin (Bin (Leaf 3) (Leaf 7)) (Leaf 11))
3

```

If a tree consists of a leaf, then *sizeBinTree* ignores the value at the leaf and returns 1 as the size of the tree. If a tree consists of two subtrees, then *sizeBinTree* returns the sum of the sizes of those subtrees as the size of the tree. Functions for computing the sum and the product of the integers at the leaves of a binary tree can be defined in a similar way. It suffices to define appropriate semantic actions *bin* and *leaf* on a type *t* (in this case *Int*) that correspond to the syntactic constructs *Bin* and *Leaf* of the datatype *BinTree*.

5.2.2. Trees for matching parentheses

Section 3.3.1 defines the datatype *Parentheses* for matching parentheses.

```

data Parentheses = Match Parentheses Parentheses
                 | Empty

```

For example, the sentence `()()` of the concrete syntax for matching parentheses is represented by the value *Match Empty (Match Empty Empty)* in the abstract syntax *Parentheses*. Remember that the abstract syntax ignores the terminal bracket symbols of the concrete syntax.

We can now define, in the same way as we did for lists and binary trees, an algebra type *ParenthesesAlgebra* and a fold function *foldParentheses*, which can be used to compute the depth (*depthParentheses*) and the width (*widthParentheses*) of matching parentheses in a compositional way. The depth of a string of matching parentheses

5. Compositionality

s is the largest number of unmatched parentheses that occurs in a substring of s . For example, the depth of the string $((()))()$ is 3. The width of a string of matching parentheses s is the the number of substrings that are matching parentheses themselves, which are not a substring of a surrounding string of matching parentheses. For example, the width of the string $((()))()$ is 2. Compositional functions on datatypes that describe the abstract syntax of a language are called syntax driven.

```

type ParenthesesAlgebra m = (m → m → m, m)
foldParentheses :: ParenthesesAlgebra m → Parentheses → m
foldParentheses (match, empty) = fold
  where
    fold (Match l r) = match (fold l) (fold r)
    fold Empty      = empty

depthParenthesesAlgebra :: ParenthesesAlgebra Int
depthParenthesesAlgebra = (λx y → max (1 + x) y, 0)

widthParenthesesAlgebra :: ParenthesesAlgebra Int
widthParenthesesAlgebra = (λ_ y → 1 + y          , 0)

depthParentheses, widthParentheses :: Parentheses → Int
depthParentheses = foldParentheses depthParenthesesAlgebra
widthParentheses = foldParentheses widthParenthesesAlgebra

parenthesesExample = Match (Match (Match Empty Empty) Empty)
                        (Match Empty
                          (Match (Match Empty Empty)
                                Empty)
                          )
                        )

? depthParentheses parenthesesExample
3
? widthParentheses parenthesesExample
3

```

Our example reveals that abstract syntax is not very well suited for interpretation by human beings. What is the concrete representation of the matching parenthesis example represented by *parenthesesExample*? It happens to be $((()))()((()))$. Fortunately, we can easily write a program that computes the concrete representation from the abstract one. We know exactly which terminals we have deleted when going from the concrete syntax to the abstract one. The algebra used by the function *a2cParentheses* simply reinserts those terminals that we have deleted. Note that *a2cParentheses* does not deal with layout such as blanks, indentation and newlines. For a simple example layout does not really matter. For large examples layout is very important: it can be used to let concrete representations look pretty.

```

a2cParenthesesAlgebra :: ParenthesesAlgebra String
a2cParenthesesAlgebra = (\xs ys → "(" ++ xs ++ ")" ++ ys, "")
a2cParentheses :: Parentheses → String
a2cParentheses = foldParentheses a2cParenthesesAlgebra

```

```

? a2cParentheses parenthesesExample
((( ))) () ( )

```

This example illustrates that a computer can easily interpret abstract syntax (something human beings have difficulties with). Strangely enough, human beings can easily interpret concrete syntax (something computers have difficulties with). What we would really like is that computers can interpret concrete syntax as well. This is the place where parsing enters the picture: computing an abstract representation from a concrete one is precisely what parsers are used for.

Consider the functions *parens* and *nesting* of Section 3.3.1 again.

```

open = symbol '('
close = symbol ')'

parens :: Parser Char Parentheses
parens =
  <$> open <*> parens <*> close <*> parens
  <|> succeed Empty

nesting :: Parser Char Int
nesting =
  <$> open <*> nesting <*> close <*> nesting
  <|> succeed 0

```

Function *nesting* could have been defined by means of function *parens* and a fold:

```

nesting' :: Parser Char Int
nesting' = depthParentheses <$> parens

```

(Remember that *depthParentheses* has been defined as a fold.) Functions *nesting* and *nesting'* compute exactly the same result. The function *nesting* is the *fusion* of the fold function with the parser *parens* from the function *nesting'*. Using laws for parsers and folds (not shown here) we can prove that the two functions are equal.

Note that function *nesting'* first builds a tree by means of function *parens*, and then flattens it by means of the fold. Function *nesting* never builds a tree, and is thus preferable for reasons of efficiency. On the other hand: in function *nesting'* we reuse the parser *parens* and the function *depthParentheses*, in function *nesting* we have to write our own parser, and convince ourselves that it is correct. So for reasons of ‘programming efficiency’ function *nesting'* is preferable. To obtain the best of both

5. Compositionality

deforestation

worlds, we would like to write function *nesting'* and have our compiler figure out that it is better to use function *nesting* in computations. The automatic transformation of function *nesting'* into function *nesting* is called *deforestation* (trees are removed). Some (very few) compilers are clever enough to perform this transformation automatically.

5.2.3. Expression trees

The matching parentheses grammar has only one nonterminal. Therefore its abstract syntax is described by a single datatype. In this section we look again at the expression grammar of Section 3.5:

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

This grammar has three nonterminals, E , T , and F . Using the approach from Section 2.6 we transform the nonterminals to datatypes:

$$\begin{aligned} \mathbf{data} \ E &= E_1 \ T \mid E_2 \ E \ T \\ \mathbf{data} \ T &= T_1 \ F \mid T_2 \ T \ F \\ \mathbf{data} \ F &= F_1 \ E \mid F_2 \ Int \end{aligned}$$

where we have translated *Digs* by the type *Int*. Note that this is a rather inconvenient and clumsy abstract syntax for expressions; the following abstract syntax is more convenient.

$$\mathbf{data} \ Expr = Con \ Int \mid Add \ Expr \ Expr \mid Mul \ Expr \ Expr$$

However, to illustrate the concept of mutual recursive datatypes, we will study the datatypes E , T , and F defined above. Since E uses T , T uses F , and F uses E , these three types are *mutually recursive*. The main datatype of the three datatypes is the one corresponding to the start symbol E . Since the datatypes are mutually recursive, the algebra type *EAlgebra* consists of three tuples of functions and three carriers (the main carrier is, as always, the one corresponding to the main datatype and is therefore the one corresponding to the start-symbol).

$$\mathbf{type} \ EAlgebra \ e \ t \ f = ((t \rightarrow e, e \rightarrow t \rightarrow e) \\ , (f \rightarrow t, t \rightarrow f \rightarrow t) \\ , (e \rightarrow f, Int \rightarrow f) \\)$$

The fold function $foldE$ for E also folds over T and F , so it uses three mutually recursive local functions.

```

foldE :: EAlgebra e t f → E → e
foldE ((e1, e2), (t1, t2), (f1, f2)) = fold
  where
    fold (E1 t) = e1 (foldT t)
    fold (E2 e t) = e2 (fold e) (foldT t)
    foldT (T1 f) = t1 (foldF f)
    foldT (T2 t f) = t2 (foldT t) (foldF f)
    foldF (F1 e) = f1 (fold e)
    foldF (F2 n) = f2 n

```

We can now use $foldE$ to write a syntax driven expression evaluator $evalE$. In the algebra that is used in the $foldE$, all type variables e , f , and t are instantiated with Int .

```

evalE :: E → Int
evalE = foldE ((id, (+)), (id, (*)), (id, id))
exE = E2 (E1 (T2 (T1 (F2 2)) (F2 3))) (T1 (F2 1))

```

```

? evalE exE
7

```

Once again our example shows that abstract syntax cannot easily be interpreted by human beings. Here is a function $a2cE$ which does this job for us.

```

a2cE :: E → String
a2cE = foldE ((e1, e2), (t1, t2), (f1, f2))
  where e1 = λt → t
        e2 = λe t → e ++ "+" ++ t
        t1 = λf → f
        t2 = λt f → t ++ "*" ++ f
        f1 = λe → "(" ++ e ++ ")"
        f2 = λn → show n

```

```

? a2cE exE
"2*3+1"

```

5. Compositionality

5.2.4. General trees

A general tree consists of a node, holding a value, where the tree splits into a list of subtrees. Notice that this list may be empty (in which case, of course, only the value at the node is of interest). As usual, the type *TreeAlgebra* and the function *foldTree* can be generated automatically from the data definition.

```
data Tree x = Node x [Tree x]
type TreeAlgebra x a = x → [a] → a
foldTree :: TreeAlgebra x a → Tree x → a
foldTree node = fold
where
    fold (Node x gts) = node x (map fold gts)
```

Notice that the constructor *Node* has a list of recursive arguments. Therefore the node function of the algebra has a corresponding list of recursive arguments. The local *fold* function is recursively called on all elements of a list using the *map* function.

One can compute the sum of the values at the nodes of a general tree as follows:

```
sumTree :: Tree Int → Int
sumTree = foldTree (λx xs → x + sum xs)
```

Computing the product of the values at the nodes of a general tree and computing the size of a general tree can be done in a similar way.

5.2.5. Efficiency

A fold takes a value of a datatype, and replaces its constructors by functions. If the evaluation of each of these functions on their arguments takes constant time, evaluation of the fold takes time linear in the number of constructors in its argument. However, some functions require more than constant evaluation time. For example, list concatenation is linear in its left argument, and it follows that if we define the function *reverse* by

```
reverse :: [a] → [a]
reverse = foldL (λx xs → xs ++ [x], [])
```

then function *reverse* takes time quadratic in the length of its argument list. So, folds are often efficient functions, but if the functions in the algebra are not constant, the fold is usually not linear. Often such a nonlinear fold can be transformed into a more efficient function. A technique that often can be used in such a case is the accumulating parameter technique. For example, for the *reverse* function we have

```
reverse x = reverse' x []
```

```

reverse' :: [a] -> [a] -> [a]
reverse' []      ys = ys
reverse' (x : xs) ys = reverse' xs (x : ys)

```

The evaluation of $reverse' xs$ takes time linear in the length of xs .

Exercise 5.1. Define an algebra type and a fold function for the following datatype.

```

data LNTree a b = Leaf a
                | Node (LNTree a b) b (LNTree a b)

```

Exercise 5.2. Define the following functions as folds on the datatype *BinTree*, see Section 5.2.1.

1. *height*, which returns the height of a tree.
2. *flatten*, which returns the list of leaf values in left-to-right order.
3. *maxBinTree*, which returns the maximal value at the leaves.
4. *sp*, which returns the length of a shortest path.
5. *mapBinTree*, which maps a function over the elements at the leaves.

Exercise 5.3. A *path* through a binary tree describes the route from the root of the tree to some leaf. We choose to represent paths by sequences of *Direction*'s:

```

data Direction = L | R

```

in such a way that taking the left subtree in an internal node will be encoded by *L* and taking the right subtree will be encoded by *R*. Define a compositional function *allPaths* which produces all paths of a given tree. Define this function first using explicit recursion, and then using a fold.

Exercise 5.4. This exercise deals with resistors. There are some basic resistors with a fixed (floating point) resistance and, given two resistors, they can be put in parallel (*[:]*) or in sequence (*[:*]*).

1. Define the datatype *Resist* to represent resistors. Also, define the type *ResistAlgebra* and the corresponding function *foldResist*.
2. Define a compositional function *result* which determines the resistance of a resistors. (Recall the rules $\frac{1}{r} = \frac{1}{r_1} + \frac{1}{r_2}$ and $r = r_1 + r_2$.)

5.3. Algebraic semantics

This section summarizes what we have seen before, and establishes terminology.

In the previous section, we have seen how to associate an algebra type and a fold function with every datatype, or family of mutually recursive datatypes. The algebra is a tuple (one component per datatype) of tuples (one component per constructor

5. Compositionality

carrier

of the datatype) of semantic actions. The algebra is parameterized over the result type of the computation, also called the *carrier* of the algebra. If we work with a family of mutually recursive datatypes, we need one carrier per datatype, and the algebra is parameterized over all of them. We call the carrier corresponding to the main datatype of the family the *main carrier*, and the others *auxiliary carriers*.

compositional

The fold function traverses a value and recursively replaces syntactic constructors of the datatypes by the corresponding semantic actions of the algebra. Functions which are defined in terms of a fold function and an algebra are called *compositional functions*.

initial algebra

There is one special algebra: the one whose components are the constructor functions of the datatypes we operate on. This algebra, when applied via a fold, defines the identity function, and is called the *initial algebra*.

algebra
homomorphism

The function resulting from applying an algebra to a fold function is sometimes also called an *algebra homomorphism*, and the compositional function is said to define *algebraic semantics*.

algebraic semantics

5.4. Expressions

The first part of this section presents a basic expression evaluator. The evaluator is extended with variables in the second part and with local definitions in the third part.

5.4.1. Evaluating expressions

In this subsection we start with a more involved example: an expression evaluator. We will use another datatype for expressions than the one introduced in Section 5.2.3: here we will use a single, and hence non mutual recursive datatype for expressions. We restrict ourselves to float valued expressions on which addition, subtraction, multiplication and division are defined. The datatype and the corresponding algebra type and fold function are as follows:

```
infixl 7 'Mul'  
infix 7 'Dvd'  
infixl 6 'Add', 'Min'  
data Expr = Expr 'Add' Expr  
          | Expr 'Min' Expr  
          | Expr 'Mul' Expr  
          | Expr 'Dvd' Expr  
          | Num Float  
type ExprAlgebra a = (a → a → a — add
```

```

, a → a → a — min
, a → a → a — mul
, a → a → a — dvd
, Float → a) — num

```

```

foldExpr :: ExprAlgebra a → Expr → a
foldExpr (add, min, mul, dvd, num) = fold

```

where

```

fold (expr1 'Add' expr2) = fold expr1 'add' fold expr2
fold (expr1 'Min' expr2) = fold expr1 'min' fold expr2
fold (expr1 'Mul' expr2) = fold expr1 'mul' fold expr2
fold (expr1 'Dvd' expr2) = fold expr1 'dvd' fold expr2
fold (Num n)                = num n

```

There is nothing special to notice about these definitions except, perhaps, the fact that *Expr* does not have an extra parameter *x* like the list and tree examples. Computing the result of an expression now simply consists of replacing the constructors by appropriate functions.

```

resultExpr :: Expr → Float
resultExpr = foldExpr ((+), (-), (*), (/), id)

```

5.4.2. Adding variables

Our next goal is to extend the evaluator of the previous subsection such that it can handle variables as well. The values of variables are typically looked up in an environment which binds the names of the variables to values. We implement an environment as a list of name-value pairs. For our purposes names are strings and values are floats. In the following programs we will use the following functions and types:

```

type Env name value = [(name, value)]
(?) :: Eq name ⇒ Env name value → name → value
env ? x = head [v | (y, v) ← env, x == y]
type Name = String
type Value = Float

```

The datatype and the corresponding algebra type and fold function are now as follows. Note that we use the same name (*Expr*) for the datatype, although it differs from the previous *Expr* datatype.

```

data Expr = Expr 'Add' Expr
          | Expr 'Min' Expr
          | Expr 'Mul' Expr

```

5. Compositionality

```

      | Expr 'Dvd' Expr
      | Num Value
      | Var Name
type ExprAlgebra a = (a → a → a — add
                      , a → a → a — min
                      , a → a → a — mul
                      , a → a → a — dvd
                      , Value → a — num
                      , Name → a) — var

foldExpr :: ExprAlgebra a → Expr → a
foldExpr (add, min, mul, dvd, num, var) = fold
  where
    fold (expr1 'Add' expr2) = fold expr1 'add' fold expr2
    fold (expr1 'Min' expr2) = fold expr1 'min' fold expr2
    fold (expr1 'Mul' expr2) = fold expr1 'mul' fold expr2
    fold (expr1 'Dvd' expr2) = fold expr1 'dvd' fold expr2
    fold (Num n)                = num n
    fold (Var x)                 = var x

```

The datatype *Expr* now has an extra constructor: the unary constructor *Var*. Similarly, the argument of *foldExpr* now has an extra component: the unary function *var* which corresponds to the unary constructor *Var*. Computing the result of an expression somehow needs to use an environment. Here is a first, bad way of doing this: one can use it as an argument of a function that computes an algebra (we will explain why this is a bad choice in the next subsection; the basic idea is that we use the environment as a global variable here).

```

resultExprBad :: Env Name Value → Expr → Value
resultExprBad env = foldExpr ((+), (-), (*), (/), id, (env?))

```

```

? resultExprBad [("x",3)] (Var "x" 'Mul' Num 2)
6

```

The good way of using an environment is the following: instead of working with a computation which, given an environment, yields an algebra of values it is better to turn the computation itself into an algebra. Thus we turn the environment in a 'local' variable.

```

(|+|), (|-|), (|*|), (|/|) :: (Env Name Value → Value) →
                             (Env Name Value → Value) →
                             (Env Name Value → Value)
f |+| g = λenv → f env + g env
f |-| g = λenv → f env - g env

```

```

f |*| g = λenv → f env * g env
f |/| g = λenv → f env / g env
resultExprGood :: Expr → (Env Name Value → Value)
resultExprGood = foldExpr ((|+|), (|-|), (|*|), (|/|), const, flip (?))

```

```

? resultExprGood (Var "x" 'Mul' Num 2) [("x",3)]
6

```

The actions (+), (−), (∗) and (/) on values are now replaced by corresponding actions (|+|), (|-|), (|*|) and (|/|) on computations. Computing the result of the sum of two subexpressions within a given environment consists of computing the result of the subexpressions within this environment and adding both results to yield a final result. Computing the result of a constant does not need the environment at all. Computing the result of a variable consists of looking it up in the environment. Thus, the algebraic semantics of an expression is a *computation* which yields a value.

In this case the computation is of the form $env \rightarrow val$. The value type is an example of a *synthesised attribute*. The value of an expression is synthesised from values of its subexpressions. The environment type is an example of an *inherited attribute*. The environment which is used by the computation of a subexpression of an expression is inherited from the computation of the expression. Since we are working with abstract syntax we say that the synthesised and inherited attributes are attributes of the datatype *Expr*. If *Expr* is one of the mutually recursive datatypes which are generated from the nonterminals of a grammar, then we say that the synthesised and inherited attributes are attributes of the nonterminal.

synthesised attribute
inherited attribute

5.4.3. Adding definitions

Our next goal is to extend the evaluator of the previous subsection such that it can handle definitions as well. A (local) definition is an expression of the form

$$Def \ name \ expr_1 \ expr_2$$

which should be interpreted as: let the value of *name* be equal to $expr_1$ in expression $expr_2$. Variables are typically defined by updating the environment with an appropriate name-value pair.

The datatype (called *Expr* again) and the corresponding algebra type and eval function are now as follows:

```

data Expr = Expr 'Add' Expr
          | Expr 'Min' Expr
          | Expr 'Mul' Expr

```

5. Compositionality

```

| Expr 'Dvd' Expr
| Num Value
| Var Name
| Def Name Expr Expr
type ExprAlgebra a = (a → a → a           — add
                    , a → a → a           — min
                    , a → a → a           — mul
                    , a → a → a           — dvd
                    , Value → a           — num
                    , Name → a            — var
                    , Name → a → a → a) — def

foldExpr :: ExprAlgebra a → Expr → a
foldExpr (add, min, mul, dvd, num, var, def) = fold
where
  fold (expr1 'Add' expr2) = fold expr1 'add' fold expr2
  fold (expr1 'Min' expr2) = fold expr1 'min' fold expr2
  fold (expr1 'Mul' expr2) = fold expr1 'mul' fold expr2
  fold (expr1 'Dvd' expr2) = fold expr1 'dvd' fold expr2
  fold (Num n)                = num n
  fold (Var x)                 = var x
  fold (Def x value body)     = def x (fold value) (fold body)

```

Expr now has an extra constructor: the ternary constructor *Def*, which can be used to introduce a local variable. For example, the following expression can be used to compute the number of seconds per year.

```

seconds = Def "days_per_year" (Num 365) (
  Def "hours_per_day" (Num 24) (
    Def "minutes_per_hour" (Num 60) (
      Def "seconds_per_minute" (Num 60) (
        Var "days_per_year" 'Mul'
        Var "hours_per_day" 'Mul'
        Var "minutes_per_hour" 'Mul'
        Var "seconds_per_minute"
      ))))

```

Similarly, the parameter of *foldExpr* now has an extra component: the ternary function *def* which corresponds to the ternary constructor *Def*. Notice that the last two arguments are recursive ones. We can now explain why the first use of environments is inferior to the second one. Trying to extend the first definition gives something like:

```

resultExprBad :: Env Name Value → Expr → Value
resultExprBad env =
  foldExpr ((+), (-), (*), (/), id, (env?), error "def")

```


The last component causes a problem: a body that contains a local definition has to be evaluated in an updated environment. We cannot update the environment in this setting: we can read the environment but afterwards it is not accessible any more in the algebra (which consists of values). Extending the second definition causes no problems: the environment is now accessible in the algebra (which consists of computations). We can easily add a new action which updates the environment. The computation corresponding to the body of an expression with a local definition can now be evaluated in the updated environment.

```

f |+| g = λenv → f env + g env
f |-| g = λenv → f env - g env
f |*| g = λenv → f env * g env
f | / | g = λenv → f env / g env
x |=| f = λg env → g ((x, f env) : env)
resultExprGood :: Expr → (Env Name Value → Value)
resultExprGood =
  foldExpr ((+|), (|-|), (|*|), (| / |), const, flip (?), (|=|))

```

```

? resultExprGood seconds []
31536000

```

Note that by prepending a pair (x, y) to an environment (in the definition of the operator $|=|$), we add the pair to the environment. By definition of $(?)$, the binding for x hides possible other bindings for x .

5.4.4. Compiling to a stack machine

In this section we compile expressions to instructions on a stack machine. We can then use this stack machine to evaluate compiled expressions. This section is inspired by an example in the textbook by Bird and Wadler [3].

Imagine a simple computer for evaluating arithmetic expressions. This computer has a ‘stack’ and can execute ‘instructions’ which change the value of the stack. The class of possible instructions is defined by the following datatype.

```

data MachInstr v = Push v | Apply (v → v → v)
type MachProgr v = [MachInstr v]

```

An instruction either pushes a value of type v on the stack, or it executes an operator that takes the two top values of the stack, applies the operator, and pushes the result back on the stack. A stack (a value of type $Stack\ v$ for some value type v) is a list of values, from which you can *pop* values, on which you can *push* values, and from which

5. Compositionality

you can take the *top* value. A module *Stack* for stacks can be found in Appendix A. The effect of executing an instruction of type *MachInstr* is defined by

$$\begin{aligned}
 \text{execute} &:: \text{MachInstr } v \rightarrow \text{Stack } v \rightarrow \text{Stack } v \\
 \text{execute } (\text{Push } x) & \quad s = \text{push } x \ s \\
 \text{execute } (\text{Apply } op) & \quad s = \mathbf{let} \ a = \text{top } s \\
 & \quad \quad \quad \quad \quad \quad \quad \quad t = \text{pop } s \\
 & \quad \quad \quad \quad \quad \quad \quad \quad b = \text{top } t \\
 & \quad \quad \quad \quad \quad \quad \quad \quad u = \text{pop } t \\
 & \quad \quad \quad \quad \quad \quad \quad \quad \mathbf{in} \ \text{push } (op \ a \ b) \ u
 \end{aligned}$$

A sequence of instructions is executed by the function *run* defined by

$$\begin{aligned}
 \text{run} &:: \text{MachProgr } v \rightarrow \text{Stack } v \rightarrow \text{Stack } v \\
 \text{run } [] & \quad s = s \\
 \text{run } (x : xs) & \quad s = \text{run } xs \ (\text{execute } x \ s)
 \end{aligned}$$

It follows that *run* can be defined as a *foldl*.

An expression can be translated (or compiled) into a list of instructions by the function *compile*, defined by:

$$\begin{aligned}
 \text{compileExpr} &:: \text{Expr} \rightarrow \\
 & \quad \text{Env Name (MachProgr Value)} \rightarrow \\
 & \quad \text{MachProgr Value} \\
 \text{compileExpr} &= \text{foldExpr } (\text{add}, \text{min}, \text{mul}, \text{dvd}, \text{num}, \text{var}, \text{def}) \\
 \mathbf{where} \\
 f \text{ 'add' } g &= \lambda env \rightarrow f \ env \ ++ \ g \ env \ ++ \ [\text{Apply } (+)] \\
 f \text{ 'min' } g &= \lambda env \rightarrow f \ env \ ++ \ g \ env \ ++ \ [\text{Apply } (-)] \\
 f \text{ 'mul' } g &= \lambda env \rightarrow f \ env \ ++ \ g \ env \ ++ \ [\text{Apply } (*)] \\
 f \text{ 'dvd' } g &= \lambda env \rightarrow f \ env \ ++ \ g \ env \ ++ \ [\text{Apply } (/)] \\
 \text{num } v &= \lambda env \rightarrow [\text{Push } v] \\
 \text{var } x &= \lambda env \rightarrow env \ ? \ x \\
 \text{def } x \ fd \ fb &= \lambda env \rightarrow fb \ ((x, fd \ env) : env)
 \end{aligned}$$

Exercise 5.5. Define the following functions as folds on the datatype *Expr* that contains definitions.

1. *isSum*, which determines whether or not an expression is a sum.
2. *vars*, which returns the list of variables that occur in the expression.

Exercise 5.6. This exercise deals with expressions without definitions. The function *der* is defined by

$$\begin{aligned}
 \text{der } (e_1 \text{ 'Add' } e_2) \ dx &= \text{der } e_1 \ dx \ \text{'Add' } \ \text{der } e_2 \ dx \\
 \text{der } (e_1 \text{ 'Min' } e_2) \ dx &= \text{der } e_1 \ dx \ \text{'Min' } \ \text{der } e_2 \ dx \\
 \text{der } (e_1 \text{ 'Mul' } e_2) \ dx &= (e_1 \ \text{'Mul' } \ \text{der } e_2 \ dx) \ \text{'Add' } \ (\text{der } e_1 \ dx \ \text{'Mul' } \ e_2)
 \end{aligned}$$

$$\begin{aligned} \text{der } (e_1 \text{ 'Dvd' } e_2) \text{ } dx &= ((e_2 \text{ 'Mul' } \text{der } e_1 \text{ } dx) \text{ 'Min' } (e_1 \text{ 'Mul' } \text{der } e_2 \text{ } dx)) \\ &\quad \text{'Dvd' } (e_2 \text{ 'Mul' } e_2) \\ \text{der } (\text{Num } f) \quad dx &= \text{Num } 0 \\ \text{der } (\text{Var } s) \quad dx &= \text{if } s == dx \text{ then Num } 1 \text{ else Num } 0 \end{aligned}$$

1. Give an informal description of the function *der*.
2. Why is the function *der* not compositional ?
3. Define a datatype *Exp* to represent expressions consisting of (floating point) constants, variables, addition and subtraction. Also, define the type *ExpAlgebra* and the corresponding *foldExp*.
4. Define the function *der* on *Exp* and show that this function is compositional.

Exercise 5.7. Define the function *replace*, which given a binary tree and an element *m* replaces the elements at the leaves by *m* as a fold on the datatype *BinTree*, see Section 5.2.1. It is easy to write a function with the required functionality if you swap the arguments, but then it is impossible to write *replace* as a fold. Note that the fold returns a function, which when given *m* replaces all the leaves by *m*.

Exercise 5.8. Consider the datatype of paths introduced in Exercise 5.3. A path in a tree leads to a unique leaf. Define a compositional function *path2Value* which, given a tree and a path in the tree, yields the element at the unique leaf.

5.5. Block structured languages

This section presents a more complex example of the use of tuples in combination with compositionality. The example deals with the scope of variables in a block structured language. A variable from a global scope is visible in a local scope only if it is not hidden by a variable with the same name in the local scope.

5.5.1. Blocks

A block is a list of statements. A statement is a variable declaration, a variable usage or a nested block. The concrete representation of an example block of our block structured language looks as follows (*dcl* stands for declaration, *use* stands for usage and *x*, *y* and *z* are variables).

```
use x ; dcl x ;
(use z ; use y ; dcl x ; dcl z ; use x) ;
dcl y ; use y
```

5. Compositionality

Statements are separated by semicolons. Nested blocks are surrounded by parentheses. The usage of z refers to the local declaration (the only declaration of z). The usage of y refers to the global declaration (the only declaration of y). The local usage of x refers to the local declaration and the global usage of x refers to the global declaration. Note that it is allowed to use variables before they are declared. Here are some mutually recursive (data)types, which describe the abstract syntax of blocks, corresponding to the grammar that describes the concrete syntax of blocks which is used above. We use meaningful names for data constructors and we use built-in lists instead of user-defined lists for the block algebra. As usual, the algebra type *BlockAlgebra*, which consists of two tuples of functions, and the fold function *foldBlock*, which uses two mutually recursive local functions, can be generated from the two mutually recursive (data)types.

```

type Block      = [Statement]
data Statement = Dcl Idf | Use Idf | Blk Block
type Idf       = String
type BlockAlgebra b s = ((s → b → b, b)
                        , (Idf → s, Idf → s, b → s)
                        )
foldBlock :: BlockAlgebra b s → Block → b
foldBlock ((cons, empty), (dcl, use, blk)) = fold
where
  fold (s : b)    = cons (foldS s) (fold b)
  fold []         = empty
  foldS (Dcl x)  = dcl x
  foldS (Use x)  = use x
  foldS (Blk b)  = blk (fold b)

```

5.5.2. Generating code

The goal of this section is to generate code from a block. The code consists of a sequence of instructions. There are three types of instructions.

- *Enter* (l, c): enters the l 'th nested block in which c local variables are declared.
- *Leave* (l, c): leaves the l 'th nested block in which c local variables were declared.
- *Access* (l, c): accesses the c 'th variable of the l 'th nested block.

The code generated for the above example looks as follows.

```

[Enter (0, 2), Access (0, 0)
, Enter (1, 2), Access (1, 1), Access (0, 1), Access (1, 0), Leave (1, 2)

```

```

, Access (0, 1), Leave (0, 2)
]

```

Note that we start numbering levels (l) and counts (c) (which are sometimes called displacements) from 0. The abstract syntax of the code to be generated is described by the following datatype.

```

type Count = Int
type Level = Int
type Variable = (Level, Count)
type BlockInfo = (Level, Count)
data Instruction = Enter BlockInfo
                  | Leave BlockInfo
                  | Access Variable
type Code = [Instruction]

```

The function *ab2ac*, which generates abstract code (a value of type *Code*) from an abstract block (a value of type *Block*), uses a compositional function *block2Code*. For all syntactic constructs of *Statement* and *Block* we define appropriate semantic actions on an algebra of computations. Here is a, somewhat simplified, description of these semantic actions.

- *Dcl*: Every time we declare a local variable x we have to update the local environment le of the block we are in by associating with x the current pair of level and local-count (l, lc) . Moreover we have to increment the local variable count lc to $lc + 1$. Note that we do not generate any code for a declaration statement. Instead we perform some computations which make it possible to generate appropriate code for other statements.

```

dcl  $x$  ( $le, l, lc$ ) = ( $le', lc'$ )
where
   $le' = le$  'update' ( $x, (l, lc)$ )
   $lc' = lc + 1$ 

```

where function *update* is defined in the *AssociationList* module.

- *Use*: Every time we use a local variable x we have to generate code cd' for it. This code is of the form [*Access* (l, lc)]. The level–local-count pair (l, lc) of the variable is looked up in the global environment e .

```

use  $x$   $e = cd'$ 
where
   $cd' = [Access (l, c)]$ 
   $(l, c) = e ? x$ 

```

5. Compositionality

- *Blk*: Every time we enter a nested block we increment the global level l to $l+1$, start with a fresh local variable count 0 and set the local environment of the nested block we enter to the current global environment e . The computation for the nested block results in a local variable count lcB and a local environment leB . Furthermore we need to make sure that the global environment (the one in which we look up variables) which is used by the computation for the nested block is equal to leB . The code which is generated for the block is surrounded by an appropriate $[Enter\ lcB]$ - $[Leave\ lcB]$ pair.

$$\begin{aligned}
 blk\ fB\ (e, l) &= cd' \\
 \text{where} \\
 l' &= l + 1 \\
 (leB, lcB, cdB) &= fB\ (leB, l', e, 0) \\
 cd' &= [Enter\ (l', lcB)] \text{++} cdB \text{++} [Leave\ (l', lcB)]
 \end{aligned}$$

- $[]$: No action need to be performed for an empty block.
- $(:)$: For every nonempty block we perform the computation of the first statement of the block which, given a local environment le and local variable count lc , results in a local environment leS and local variable count lcS . This environment-count pair is then used by the computation of the rest of the block to result in a local environment le' and local variable count lc' . The code cd' which is generated is the concatenation $cdS \text{++} cdB$ of the code cdS which is generated for the first statement and the code cdB which is generated for the rest of the block.

$$\begin{aligned}
 cons\ fS\ fB\ (le, lc) &= (le', lc', cd') \\
 \text{where} \\
 (leS, lcS, cdS) &= fS\ (le, lc) \\
 (le', lc', cdB) &= fB\ (leS, lcS) \\
 cd' &= cdS \text{++} cdB
 \end{aligned}$$

What does our actual computation type look like? For *dcl* we need three inherited attributes: a global level, a local block environment and a local variable count. Two of them: the local block environment and the local variable count are also synthesised attributes. For *use* we need one inherited attribute: a global block environment, and we compute one synthesised attribute: the generated code. For *blk* we need two inherited attributes: a global block environment and a global level, and we compute two synthesised attributes: the local variable count and the generated code. Moreover there is one extra attribute: a local block environment which is both inherited and synthesised. When processing the statements of a nested block we already make use of the global block environment which we are synthesising (when looking up variables). For *cons* we compute three synthesised attributes: the local block environment, the local variable count and the generated code. Two of them, the

local block environment and the local variable count are also needed as inherited attributes. It is clear from the considerations above that the following types fulfill our needs.

```

type BlockEnv = [(Idf, Variable)]
type GlobalEnv = (BlockEnv, Level)
type LocalEnv = (BlockEnv, Count)

```

The implementation of *block2Code* is now a straightforward translation of the actions described above. Attributes which are not mentioned in those actions are added as extra components which do not contribute to the functionality.

```

block2Code :: Block → GlobalEnv → LocalEnv → (LocalEnv, Code)
block2Code = foldBlock ((cons, empty), (dcl, use, blk))
where
  cons fS fB (e, l) (le, lc) = ((le', lc'), cd')
    where
      ((leS, lcS), cdS) = fS (e, l) (le, lc)
      ((le', lc'), cdB) = fB (e, l) (leS, lcS)
      cd' = cdS ++ cdB
  empty (e, l) (le, lc) = ((le, lc), [])
  dcl x (e, l) (le, lc) = ((le', lc'), [])
    where
      le' = (x, (l, lc)) : le
      lc' = lc + 1
  use x (e, l) (le, lc) = ((le, lc), cd')
    where
      cd' = [Access (l, c)]
      (l, c) = e ? x
  blk fB (e, l) (le, lc) = ((le, lc), cd')
    where
      ((leB, lcB), cdB) = fB (leB, l') (e, 0)
      l' = l + 1
      cd' = [Enter (l', lcB)] ++ cdB ++ [Leave (l', lcB)]

```

The code generator starts with an empty local environment, a fresh level and a fresh local variable count. The code is a synthesised attribute. The global environment is an attribute which is both inherited and synthesised. When processing a block we already use the global environment which we are synthesising (when looking up variables).

```

ab2ac :: Block → Code
ab2ac b = [Enter (0, c)] ++ cd ++ [Leave (0, c)]
where

```

5. Compositionality

```
((e, c), cd) = block2Code b (e, 0) ([], 0)
aBlock
= [ Use "x", Dcl "x"
  , Blk [Use "z", Use "y", Dcl "x", Dcl "z", Use "x"]
  , Dcl "y", Use "y"]

? ab2ac aBlock
[Enter (0,2), Access (0,0)
, Enter (1,2), Access (1,1), Access (0,1), Access (1,0), Leave (1,2)
, Access (0,1), Leave (0,2)]
```

5.6. Exercises

Exercise 5.9. Consider your answer to Exercise 2.22, which gives an abstract syntax for palindromes.

1. Define a type *PalAlgebra* that describes the type of the semantic actions that correspond to the syntactic constructs of *Pal*.
2. Define the function *foldPal*, which describes how the semantics actions that correspond to the syntactic constructs of *Pal* should be applied.
3. Define the functions *a2cPal* and *aCountPal* as *foldPal*'s.
4. Define the parser *pfoldPal* which interprets its input in an arbitrary semantic *PalAlgebra* without building the intermediate abstract syntax tree.
5. Describe the parsers *pfoldPal* m_1 and *pfoldPal* m_2 where m_1 and m_2 correspond to the algebras of *a2cPal* and *aCountPal* respectively.

Exercise 5.10. Consider your answer to Exercise 2.23, which gives an abstract syntax for mirror-palindromes.

1. Define the type *MirAlgebra* that describes the semantic actions that correspond to the syntactic constructs of *Mir*.
2. Define the function *foldMir*, which describes how semantic actions that correspond to the syntactic constructs of *Mir* should be applied.
3. Define the functions *a2cMir* and *m2pMir* as *foldMir*'s.
4. Define the parser *pfoldMir*, which interprets its input in an arbitrary semantic *MirAlgebra* without building the intermediate abstract syntax tree.
5. Describe the parsers *pfoldMir* m_1 and *pfoldMir* m_2 where m_1 and m_2 correspond to the algebras of *a2cMir* and *m2pMir*, respectively.

Exercise 5.11. Consider your answer to exercise 2.24, which gives an abstract syntax for parity-sequences.

1. Define the type *ParityAlgebra* that describes the semantic actions that correspond to the syntactic constructs of *Parity*.
2. Define the function *foldParity*, which describes how the semantic actions that correspond to the syntactic constructs of *Parity* should be applied.
3. Define the function *a2cParity* as *foldParity*.

Exercise 5.12. Consider your answer to Exercise 2.25, which gives an abstract syntax for bit-lists.

1. Define the type *BitListAlgebra* that describes the semantic actions that correspond to the syntactic constructs of *BitList*.
2. Define the function *foldBitList*, which describes how the semantic actions that correspond to the syntactic constructs of *BitList* should be applied.
3. Define the function *a2cBitList* as a *foldBitList*.
4. Define the parser *pfoldBitList*, which interprets its input in an arbitrary semantic *BitListAlgebra* without building the intermediate abstract syntax tree.

Exercise 5.13. The following grammar describes the concrete syntax of a simple block-structured programming language

$$\begin{array}{ll}
 B \rightarrow S R & (\text{block}) \\
 R \rightarrow ; S R \mid \varepsilon & (\text{rest}) \\
 S \rightarrow D \mid U \mid N & (\text{statement}) \\
 D \rightarrow \mathbf{x} \mid \mathbf{y} & (\text{declaration}) \\
 U \rightarrow \mathbf{X} \mid \mathbf{Y} & (\text{usage}) \\
 N \rightarrow (B) & (\text{nested block})
 \end{array}$$

1. Define a datatype *Block* that describes the abstract syntax that corresponds to the grammar. What is the abstract representation of $\mathbf{x}; (\mathbf{y}; \mathbf{Y}); \mathbf{X}$?
2. Define the type *BlockAlgebra* that describes the semantic actions that correspond to the syntactic constructs of *Block*.
3. Define the function *foldBlock*, which describes how the semantic actions corresponding to the syntactic constructs of *Block* should be applied.
4. Define the function *a2cBlock*, which converts an abstract block into a concrete one. Write *a2cBlock* as a *foldBlock*.
5. The function *checkBlock* tests whether or not each variable of a given abstract block is declared before use (declared in the same or in a surrounding block).

5. Compositionality

6. Computing with parsers

Parsers produce results. For example, the parsers for travelling schemes given in Chapter 4 return an abstract syntax, or an integer that represents the net travelling time in minutes. The net travelling time is computed directly by inserting the correct semantic functions. Another way to compute the net travelling time is by first computing the abstract syntax, and then applying a function to the abstract syntax that computes the net travelling time. This section shows several ways to compute results using parsers:

- insert a semantic function in the parser;
- apply a fold to the abstract syntax;
- use a class instead of abstract syntax;
- pass an algebra to the parser.

6.1. Insert a semantic function in the parser

In Chapter 4 we have defined two parsers: a parser that computes the abstract syntax for a travelling schema, and a parser that computes the net travelling time. These functions are obtained by inserting different functions in the basic parser. If we want to compute the total travelling time, we have to insert different functions in the basic parser. This approach works fine for a small parser, but it has some disadvantages when building a larger parser:

- semantics is intertwined with the parsing process;
- it is difficult to locate all positions where semantic functions have to be inserted in the parser.

6.2. Apply a fold to the abstract syntax

Instead of inserting operations in a basic parser, we can write a parser that parses the input to an abstract syntax, and computes the desired result by applying a fold to the abstract syntax.

An example of such an approach has been given in Section 5.2.2, where we defined two functions with the same functionality: `nesting` and `nesting'`; both compute the maximum nesting depth in a string of parentheses. Function `nesting` is defined

6. Computing with parsers

by inserting functions in the basic parser. Function `nesting'` is defined by applying a fold to the abstract syntax. Each of these definitions has its own merits; we repeat the main arguments below.

```
parens  :: Parser Char Parentheses
parens  = (\_ b _ d -> Match b d) <$>
          open <*> parens <*> close <*> parens
          <|> succeed Empty

nesting :: Parser Char Int
nesting = (\_ b _ d -> max (1+b) d) <$>
          open <*> nesting <*> close <*> nesting
          <|> succeed 0

nesting' :: Parser Char Int
nesting' = depthParentheses <$> parens
```

The first definition (`nesting`) is more efficient, because it does not build an intermediate abstract syntax tree. On the other hand, it might be more difficult to write because we have to insert functions in the correct places in the basic parser. The advantage of the second definition (`nesting'`) is that we reuse both the parser `parens`, which returns an abstract syntax tree, and the function `depthParentheses` (or the function `foldParentheses`, which is used in the definition of `depthParentheses`), which does recursion over an abstract syntax tree. The only thing we have to write ourselves in the second definition is the `depthParenthesesAlgebra`. The disadvantage of the second definition is that it builds an intermediate abstract syntax tree, which is ‘flattened’ by the fold. We want to avoid building the abstract syntax tree altogether. To obtain the best of both worlds, we would like to write function `nesting'` and have our compiler figure out that it is better to use function `nesting` in computations. The automatic transformation of function `nesting'` into function `nesting` is called *deforestation* (trees are removed). Some (very few) compilers are clever enough to perform this transformation automatically.

6.3. Deforestation

Deforestation removes intermediate trees in computations. The previous section gives an example of deforestation on the datatype `Parentheses`. This section sketches the general idea.

Suppose we have a datatype `AbstractTree`

```
data AbstractTree = ...
```

From this datatype we construct an algebra and a fold, see Chapter 5.

```
type AbstractTreeAlgebra a = ...

foldAbstractTree :: AbstractTreeAlgebra a -> AbstractTree -> a
```

A parser for the datatype `AbstractTree` (which returns a value of `AbstractTree`) has the following type:

```
parseAbstractTree :: Parser Symbol AbstractTree
```

where `Symbol` is some type of input symbols (for example `Char`). Suppose now that we define a function `p` that parses an `AbstractTree`, and then computes some value by folding with an algebra `f` over this tree:

```
p = foldAbstractTree f . parseAbstractTree
```

Then deforestation says that `p` is equal to the function `parseAbstractTree` in which occurrences of the constructors of the datatype `AbstractTree` have been replaced by the corresponding components of the algebra `f`. The following two sections each describe a way to implement such a deforested function.

6.4. Using a class instead of abstract syntax

Classes can be used to implement the deforested or fused computation of a fold with a parser. This gives a solution of the desired efficiency.

For example, for the language of parentheses, we define the following class:

```
class Parens a where
  match :: a -> a -> a
  empty :: a
```

Note that types of the functions in the class `Parens` correspond exactly to the two types that occur in the type `ParenthesesAlgebra`. This class is used in a parser for parentheses:

```
parens :: Parens a => Parser Char a
parens = (\_ b _ d -> match b d) <$>
        open <*> parens <*> close <*> parens
        <|> succeed empty
```

6. Computing with parsers

The algebra is implicit in this function: the only thing we know is that there exist functions `empty` and `match` of the correct type; we know nothing about their implementation. To obtain a function `parens` that returns a value of type `Parentheses` we create the following instance of the class `Parens`.

```
instance Parens Parentheses where
  match = Match
  empty = Empty
```

Now we can write:

```
?(parens :: Parser Char Parentheses) "()()"
[(Match Empty (Match Empty Empty), "")
 ,(Match Empty Empty, "()")
 ,(Empty, "()()")]
]
```

Note that we have to supply the type of `parens` in this expression, otherwise Haskell doesn't know which instance of `Parens` to use. This is how we turn the implicit 'class' algebra into an explicit 'instance' algebra. Another instance of `Parens` can be used to compute the nesting depth of parentheses:

```
instance Parens Int where
  match b d = max (1+b) d
  empty     = 0
```

And now we can write:

```
?(parens :: Parser Char Int) "()()"
[(1, ""), (1, "()"), (0, "()()")]
```

So the answer depends on the type we want our function `parens` to have. This also immediately shows a problem of this, otherwise elegant, approach: it does not work if we want to compute two different results of the same type, because Haskell doesn't allow you to define two (or more) instances with the same type. So once we have defined the instance `Parens Int` as above, we cannot use function `parens` to compute, for example, the width (also an `Int`) of a string of parentheses.

6.5. Passing an algebra to the parser

The previous section shows how to implement a parser with an implicit algebra. Since this approach fails when we want to define different parsers with the same result type, we make the algebras explicit. Thus we obtain the following definition of `parens`:

```
parens                :: ParenthesesAlgebra a -> Parser Char a
parens (match,empty) = par where
  par = (\_ b _ d -> match b d) <$>
        open <*> par <*> close <*> par
        <|> succeed empty
```

Note that it is now easy to define different parsers with the same result type:

```
nesting, breadth :: Parser Char Int
nesting          = parens (\b d -> max (1+b) d,0)
breadth          = parens (\b d -> d+1,0)
```

6. Computing with parsers

7. Programming with higher-order folds

Introduction

In the previous chapters we have seen that algebras play an important role when describing the meaning of a recognised structure (a parse tree). For each recursive datatype `T` we have a function `foldT`, and for each constructor of the datatype we have a corresponding function as a component in the algebra. Chapter 5 introduces a language in which local declarations are permitted. Evaluating expressions in this language can be done by choosing an appropriate algebra. The domain of that algebra is a higher order (data)type (a (data)type that contains functions). Unfortunately, the resulting code comes as a surprise to many. In this chapter we will illustrate a related formalism, which will make it easier to construct such involved algebras. This related formalism is the attribute grammar formalism. We will not formally define attribute grammars, but instead illustrate the formalism with some examples, and give an informal definition.

We start with developing a somewhat unconventional way of looking at functional programs, and especially those programs that use functions that recursively descend over datatypes a lot. In our case one may think about these datatypes as abstract syntax trees. When computing a property of such a recursive object (for example, a program) we define two sets of functions: one set that describes how to recursively visit the nodes of the tree, and one set of functions (an algebra) that describes what to compute at each node when visited.

One of the most important steps in this process is deciding what the carrier type of the algebras is going to be. Once this step has been taken, these types are a guideline for further design steps. We will see that such carrier types may be functions themselves, and that deciding on the type of such functions may not always be simple. In this chapter we will present a view on recursive computations that will enable us to “design” the carrier type in an incremental way. We will do so by constructing algebras out of other algebras. In this way we define the meaning of a language in a *semantically compositional* way.

We will start with the `rep_min` example, which looks a bit artificial, and deals with a non-interesting, highly specific problem. However, it has been chosen for its simplicity, and to not distract our attention to specific, programming language related,

```

data Tree = Leaf Int
          | Bin Tree Tree deriving Show

type TreeAlgebra a = (Int -> a, a -> a -> a)

foldTree :: TreeAlgebra a -> Tree -> a
foldTree alg@(leaf, _ ) (Leaf i)   = leaf i
foldTree alg@(_     , bin) (Bin l r) = bin (foldTree alg l)
                                       (foldTree alg r)

```

Listing 7.1: `rm.start.hs`

semantic issues. The second example of this chapter demonstrates the techniques on a larger example: a small compiler for part of a programming language.

Goals

In this chapter you will learn:

- how to write ‘circular’ functional programs, or ‘higher-order folds’;
- how to combine algebras;
- (informally) the concept of an attribute grammar.

7.1. The `rep_min` problem

One of the famous examples in which the power of lazy evaluation is demonstrated is the so-called *rep_min* problem [2]. Many have wondered how this program achieves its goal, since at first sight it seems that it is impossible to compute anything with this program. We will use this problem, and a sequence of different solutions, to build up an understanding of a whole class of such programs.

In Listing 7.1 we present the datatype `Tree`, together with its associated algebra. The *carrier type* of an algebra is the type that describes the objects of the algebra. We represent it by a type parameter of the algebra type:

```

type TreeAlgebra a = (Int -> a, a -> a -> a)

```

The associated evaluation function `foldTree` systematically replaces the constructors `Leaf` and `Bin` by corresponding operations from the algebra `alg` that is passed as an argument.

```

minAlg  :: TreeAlgebra Int
minAlg  = (id, min :: Int->Int->Int)

rep_min  :: Tree -> Tree
rep_min t = foldTree repAlg t
  where m      = foldTree minAlg t
        repAlg = (const (Leaf m), Bin)

```

Listing 7.2: rm.sol1.hs

We now want to construct a function `rep_min :: Tree -> Tree` that returns a `Tree` with the same “shape” as its argument `Tree`, but with the values in its leaves replaced by the minimal value occurring in the original tree. For example,

```

?rep_min (Bin (Bin (Leaf 1) (Leaf 7)) (Leaf 11))
Bin (Bin (Leaf 1) (Leaf 1)) (Leaf 1)

```

7.1.1. A straightforward solution

A straightforward solution to the `rep_min` problem consists of a function in which `foldTree` is used twice: once for computing the minimal value of the leaf values, and once for constructing the resulting `Tree`. The function `rep_min` that solves the problem in this way is given in Listing 7.2. Notice that the variable `m` is a global variable of the `repAlg`-algebra, that is used in the tree constructing call of `foldTree`. One of the disadvantages of this solution is that in the course of the computation the pattern matching associated with the inspection of the tree nodes is performed twice for each node in the tree.

Although this solution as such is no problem, we will try to construct a solution that calls `foldTree` only once.

7.1.2. Lambda lifting

We want to obtain a program for the `rep_min` problem in which pattern matching is used only once. Program Listing 7.3 is an intermediate step towards this goal. In this program the global variable `m` has been removed and the second call of `foldTree` does not construct a `Tree` anymore, but instead *a function constructing a tree* of type `Int -> Tree`, which takes the computed minimal value as an argument. Notice how we have emphasized the fact that a function is returned through some superfluous notation: the first lambda in the function definitions constituting the algebra `repAlg` is required by the signature of the algebra, the second lambda, which could have been

7. Programming with higher-order folds

```
repAlg = ( \_      -> \m -> Leaf m
          ,\lfun rfun -> \m -> let lt = lfun m
                               rt = rfun m
                               in Bin lt rt
          )

rep_min' t = (foldTree repAlg t) (foldTree minAlg t)
```

Listing 7.3: rm.sol2.hs

```
infix 9 'tuple'

tuple :: TreeAlgebra a -> TreeAlgebra b -> TreeAlgebra (a,b)
(leaf1, bin1) 'tuple' (leaf2, bin2) = (\i  -> (leaf1 i, leaf2 i)
                                       ,\l r -> (bin1 (fst l) (fst r)
                                               ,bin2 (snd l) (snd r)
                                               )
                                       )

min_repAlg :: TreeAlgebra (Int, Int -> Tree)
min_repAlg = (minAlg 'tuple' repAlg)

rep_min'' t = r m
  where (m, r) = foldTree min_repAlg t
```

Listing 7.4: rm.sol3.hs

omitted, is there because the carrier set of the algebra contains functions of type `Int -> Tree`. This process is done routinely by functional compilers and is known as *lambda-lifting*.

7.1.3. Tupling computations

We are now ready to formulate a solution in which `foldTree` is called only once. Note that in the last solution the two calls of `foldTree` don't interfere with each other. As a consequence we may perform both the computation of the tree constructing function and the minimal value in one go, by tupling the results of the computations. The solution is given in Listing 7.4. First a function `tuple` is defined. This function takes two `TreeAlgebras` as arguments and constructs a third one, which has as its carrier tuples of the carriers of the original algebras.

7.1.4. Merging tupled functions

In the next step we transform the type of the carrier set in the previous example, $(\text{Int}, \text{Int} \rightarrow \text{Tree})$, into an equivalent type $\text{Int} \rightarrow (\text{Int}, \text{Tree})$. This transformation is not essential here, but we use it to demonstrate that if we compute a cartesian product of functions, we may transform that type into a new type in which we compute only one function, which takes as its arguments the cartesian product of all the arguments of the functions in the tuple, and returns as its result the cartesian product of the result types. In our example the computation of the minimal value may be seen as a function of type $() \rightarrow \text{Int}$. As a consequence the argument of the new type is $((), \text{Int})$, which is isomorphic to just Int , and the result type becomes $(\text{Int}, \text{Tree})$.

We want to mention here too that the reverse is in general not true; given a function of type $(a, b) \rightarrow (c, d)$, it is in general not possible to split this function into two functions of type $a \rightarrow c$ and $b \rightarrow d$, which together achieve the same effect. The new version is given in Listing 7.5.

Notice how we have, in an attempt to make the different rôles of the parameters explicit, again introduced extra lambdas in the definition of the functions of the algebra. The parameters after the second lambda are there because we construct values in a higher order carrier set. The parameters after the first lambda are there because we deal with a `TreeAlgebra`. A curious step taken here is that part of the result, in our case the value `m`, is passed back as an argument to the result of `(foldTree mergedAlg t)`. Lazy evaluation makes this work.

That such programs were possible came originally as a great surprise to many functional programmers, especially to those who used to program in LISP or ML, languages that require arguments to be evaluated completely before a call is evaluated (so-called *strict evaluation* in contrast to lazy evaluation). Because of this surprising behaviour this class of programs became known as *circular programs*. Notice however that there is nothing circular in this program. Each value is defined in terms of other values, and no value is defined in terms of itself (as in `ones=1:ones`).

Finally, Listing 7.6 shows the version of this program in which the function `foldTree` has been unfolded. Thus we obtain the original solution as given in Bird [2].

Concluding, we have systematically transformed a program that inspects each node twice into an equivalent program that inspects each node only once. The resulting solution passes back part of the result of a call as an argument to that same call. Lazy evaluation makes this possible.

Exercise 7.1. The *deepest_front* problem is the problem of finding the so-called *front* of a tree. The front of a tree is the list of all nodes that are at the deepest level. As in the `rep_min` problem, the trees involved are elements of the datatype `Tree`, see Listing 7.1. A straightforward solution is to compute the height of the tree and passing the result of this function to a function `frontAtLevel :: Tree -> Int -> [Int]`.

7. Programming with higher-order folds

```
mergedAlg :: TreeAlgebra (Int -> (Int,Tree))
mergedAlg = (\i      -> \m -> (i, Leaf m)
            ,\lfun rfun -> \m -> let (lm,lt) = lfun m
                                   (rm,rt) = rfun m
                                   in (lm 'min' rm
                                       , Bin lt rt
                                       )
            )

rep_min'''' t = r
  where (m, r) = (foldTree mergedAlg t) m
```

Listing 7.5: rm.sol4.hs

```
rep_min'''''' t = r
  where (m, r)      = tree t m
        tree (Leaf i) = \m -> (i, Leaf m)
        tree (Bin l r) = \m -> let (lm, lt) = tree l m
                                   (rm, rt) = tree r m
                                   in (lm 'min' rm, Bin lt rt)
```

Listing 7.6: rm.sol5.hs

1. Define the functions `height` and `frontAtLevel`
2. Give the four different solutions as defined in the `rep_min` problem.

Exercise 7.2. Redo the previous exercise for the *highest_front* problem.

7.2. A small compiler

This section constructs a small compiler for (a part of) a small language. The compiler compiles this code into code for a hypothetical stack machine.

7.2.1. The language

The language we consider in this section has integers, booleans, function application, and an if-then-else expression. A language with just these constructs is useless, and you will extend the language in the exercises with some other constructs, which make the language a bit more interesting. We take the following context-free grammar for the concrete syntax of the language.

$$\begin{aligned} Expr0 &\rightarrow \text{if } Expr1 \text{ then } Expr1 \text{ else } Expr1 \mid Expr1 \\ Expr1 &\rightarrow Expr2 Expr2^* \\ Expr2 &\rightarrow Int \mid Bool \end{aligned}$$

where *Int* generates integers, and *Bool* booleans. An abstract syntax for our language is given in Listing 7.7. Note that we use a single datatype for the abstract syntax instead of three datatypes (one for each nonterminal); this simplifies the code a bit. The Listing 7.7 also contains a definition of a fold and an algebra type for the abstract syntax.

A parser for expressions is given in Listing 7.8.

7.2.2. A stack machine

In section 5.4.4 we have defined a stack machine with which simple arithmetic expressions can be evaluated. Here we define a stack machine that has some more instructions. The language of the previous section will be compiled into code for this stack machine in the following section.

The stack machine we will use has the following instructions:

- it can load an integer;
- it can load a boolean;
- given an argument and a function on the stack, it can call the function on the argument;

7. Programming with higher-order folds

```
data ExprAS = If ExprAS ExprAS ExprAS
            | Apply ExprAS ExprAS
            | ConInt Int
            | ConBool Bool deriving Show

type ExprASAlgebra a = (a -> a -> a -> a
                       ,a -> a -> a
                       ,Int -> a
                       ,Bool -> a
                       )

foldExprAS :: ExprASAlgebra a -> ExprAS -> a
foldExprAS (iff,apply,conint,conbool) = fold
  where fold (If ce te ee) = iff (fold ce) (fold te) (fold ee)
        fold (Apply fe ae) = apply (fold fe) (fold ae)
        fold (ConInt i)    = conint i
        fold (ConBool b)  = conbool b
```

Listing 7.7: ExprAbstractSyntax.hs

- it can set a label in the code;
- given a boolean on the stack, it can jump to a label provided the boolean is false;
- it can jump to a label (unconditionally).

The datatype for instructions is given in Listing 7.9.

7.2.3. Compiling to the stackmachine

How do we compile the different expressions to stack machine code? We want to define a function `compile` of type

```
compile :: ExprAS -> [InstructionSM]
```

- A `ConInt i` is compiled to a `LoadInt i`.

```
compile (ConInt i) = [LoadInt i]
```
- A `ConBool b` is compiled to a `LoadBool b`.

```
compile (ConBool b) = [LoadBool b]
```

```

sptoken :: String -> Parser Char String
sptoken s = (\_ b _ -> b) <$>
    many (symbol ' ') <*> token s <*> many1 (symbol ' ')

boolean = const True <$> token "True" <|> const False <$> token "False"

parseExpr :: Parser Char ExprAS
parseExpr = expr0
  where expr0 = (\a b c d e f -> If b d f) <$>
    sptoken "if"
    <*> parseExpr
    <*> sptoken "then"
    <*> parseExpr
    <*> sptoken "else"
    <*> parseExpr
    <|> expr1
  expr1 = chain1 expr2 (const Apply <$> many1 (symbol ' '))
    <|> expr2
  expr2 = ConBool <$> boolean
    <|> ConInt <$> natural

```

Listing 7.8: ExprParser.hs

```

data InstructionSM = LoadInt Int
                  | LoadBool Bool
                  | Call
                  | SetLabel Label
                  | BrFalse Label
                  | BrAlways Label

type Label = Int

```

Listing 7.9: InstructionSM.hs

7. Programming with higher-order folds

- An application `Apply f x` is compiled by first compiling the argument `x`, then the ‘function’ `f` (at the moment it is impossible to define functions in our language, hence the quotes around ‘function’), and finally putting a `Call` on top of the stack.

```
compile (Apply f x) = compile x ++ compile f ++ [Call]
```

- An if-then-else expression `If ce te ee` is compiled by first compiling the conditional expression `ce`. Then we jump to a label (which will be set before the code of the else expression `ee` later) if the resulting boolean is false. Then we compile the then expression `te`. After the then expression we always jump to the end of the code of the if-then-else expression, for which we need another label. Then we set the label for the else expression, we compile the else expression `ee`, and, finally, we set the label for the end of the if-then-else expression.

```
compile (If ce te ee) = compile ce
                      ++ [BrFalse ?lab1]
                      ++ compile te
                      ++ [BrAlways ?lab2]
                      ++ [SetLabel ?lab1]
                      ++ compile ee
                      ++ [SetLabel ?lab2]
```

Note that we use labels here, but where do these labels come from?

From the above description we see that we also need labels when compiling an expression. We add a label argument (an integer, used for the first label in the compiled code) to function `compile`, and we want function `compile` to return the first unused label. We change the type of function `compile` as follows:

```
compile :: ExprAS -> Label -> ([InstructionSM],Label)

type Label = Int
```

The four cases in the definition of `compile` have to take care of the labels. We obtain the following definition of `compile`:

```
compile (ConInt i)      = \l -> ([LoadInt i],l)
compile (ConBool b)    = \l -> ([LoadBool b],l)
compile (Apply f x)    = \l -> let (xc,l') = compile x l
                                (fc,l'') = compile f l'
                                in (xc ++ fc ++ [Call],l'')
compile (If ce te ee)  = \l -> let (cc,l') = compile ce (l+2)
                                (tc,l'') = compile te l'
                                (ec,l''') = compile ee l''
```

```

compile = foldExprAS compileAlgebra

compileAlgebra :: ExprASAlgebra (Label -> ([InstructionSM],Label))
compileAlgebra = (\cce cte cee -> \l ->
    let (cc,l') = cce (l+2)
        (tc,l'') = cte l'
        (ec,l''') = cee l'''
    in (   cc
        ++ [BrFalse l]
        ++ tc
        ++ [BrAlways (l+1)]
        ++ [SetLabel l]
        ++ ec
        ++ [SetLabel (l+1)]
        ,l''')
    )
    ,\cf cx -> \l -> let (xc,l') = cx l
                    (fc,l'') = cf l'
                    in (xc ++ fc ++ [Call],l'')
    ,\i -> \l -> ([LoadInt i],l)
    ,\b -> \l -> ([LoadBool b],l)
    )

```

Listing 7.10: CompileExpr.hs

```

    in (   cc
        ++ [BrFalse l]
        ++ tc
        ++ [BrAlways (l+1)]
        ++ [SetLabel l]
        ++ ec
        ++ [SetLabel (l+1)]
        ,l''')
    )

```

Function `compile` is a fold, the carrier type of its algebra is a function of type `Label -> ([InstructionSM],Label)`. The definition of function `compile` as a fold is given in Listing 7.10.

Exercise 7.3 (no answer provided). Extend the code generation example by adding variables to the datatype `Expr`.

Exercise 7.4 (no answer provided). Extend the code generation example by adding definitions to the datatype `Expr` too.

7.3. Attribute grammars

In Section 7.1 we have written a program that solves the `rep_min` problem. This program computes the minimum of a tree, and it computes the tree in which all the leaf values are replaced by the minimum value. The minimum is computed bottom-up: it is *synthesized* from its children. The minimum value is then passed on to the functions that build the tree with the minimum value in its leaves. These functions receive the minimum value from their parent tree node: they *inherit* the minimum value from their parent.

We can see the `rep_min` computation as a computation on a value of type `Tree`, on which two attributes are defined: the minimum and result tree attributes. The minimum is computed bottom-up, and is then passed down to the result tree, and is therefore a synthesized and inherited attribute. The result tree is computed bottom-up, and is hence a synthesized attribute.

The formalism in which it is possible to specify such attributes and computations on datatypes or grammars is called *attribute grammars*, and was originally proposed by Donald Knuth in [9]. Attribute grammars provide a solution for the systematic description of the phases of the compiler that come after scanning and parsing. Although they look different from what we have encountered thus far and are probably a little easier to write, they can straightforwardly be mapped onto a functional program. The programs you have seen in this chapter could also have been obtained by means of such a mapping from an attribute grammar specification. Traditionally such attribute grammars are used as the input of a *compiler generator*. Just as we have seen how by introducing a suitable set of parsing combinators one may avoid the use of a special parser generator and even gain a lot of flexibility in extending the grammatical formalism by introducing more complicated combinators, we have shown how one can do without a special purpose attribute grammar processing system. But, just as the concept of a context free grammar was useful in understanding the fundamentals of parser combinators, understanding attribute grammars will help significantly in describing the semantic part of the recognition and compilation process. This chapter does not further introduce attribute grammars, but they will appear again in the course in implementing programming languages.

8. Regular Languages

Introduction

The first phase of a compiler takes an input program, and splits the input into a list of terminal symbols: keywords, identifiers, numbers, punctuation, etc. Regular expressions are used for the description of the terminal symbols. A regular grammar is a particular kind of context-free grammar that can be used to describe regular expressions. Finite-state automata can be used to recognise sentences of regular grammars. This chapter discusses all of these concepts, and is organised as follows. Section 8.1 introduces finite-state automata. Finite-state automata appear in two versions, nondeterministic and deterministic ones. Section 8.1.4 shows that a non-deterministic finite-state automaton can be transformed into a deterministic finite-state automaton, so you don't have to worry about whether or not your automaton is deterministic. Section 8.2 introduces regular grammars (context-free grammars of a particular form), and regular languages. Furthermore, it shows their equivalence with finite-state automata. Section 8.3 introduces regular expressions as finite descriptions of regular languages and shows that such expressions are another, equivalent, tool for regular languages. Finally, Section 8.4 gives some of the proofs of the results of the previous sections.

Goals

After you have studied this chapter you will know that

- regular languages are a subset of context-free languages;
- it is not always possible to give a regular grammar for a context-free grammar;
- regular grammars, finite-state automata and regular expressions are three different tools for regular languages;
- regular grammars, finite-state automata and regular expressions have the same expressive power;
- finite-state automata appear in two, equally expressive, versions: deterministic and nondeterministic.

8.1. Finite-state automata

The classical approach to recognising sentences from a regular language uses finite-state automata. A finite-state automaton can be viewed as a simple form of digital

computer with only a finite number of states, no temporary storage, an input file but only the possibility to read it, and a control unit which records the state transitions. A rather limited medium, but a useful tool in many practical subproblems. A finite-state automaton can easily be implemented by a function that takes time linear in the length of its input, and constant space. This implies that problems that can be solved by means of a finite-state automaton, can be implemented by means of very efficient programs.

8.1.1. Deterministic finite-state automata

Finite-state automata come in two flavours: deterministic and nondeterministic. We start with a description of deterministic finite-state automata, the simplest form of automata.

Definition 8.1 (Deterministic finite-state automaton, DFA). A *deterministic finite-state automaton* (DFA) is a 5-tuple (X, Q, d, S, F) where

- X is the input alphabet,
- Q is a finite set of states,
- $d :: Q \rightarrow X \rightarrow Q$ is the state transition function,
- $S \in Q$ is the start state,
- $F \subseteq Q$ is the set of accepting states.

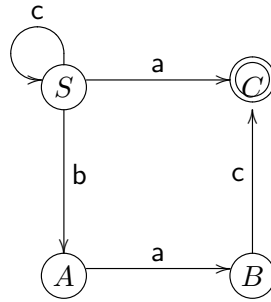
As an example, consider the DFA $M_0 = (X, Q, d, S, F)$ with

$$\begin{aligned} X &= \{a, b, c\} \\ Q &= \{S, A, B, C\} \\ F &= \{C\} \end{aligned}$$

where state transition function d is defined by

$$\begin{aligned} d S a &= C \\ d S b &= A \\ d S c &= S \\ d A a &= B \\ d B c &= C \end{aligned}$$

For human beings, a finite-state automaton is more comprehensible in a graphical representation. The following representation is customary: states are depicted as the nodes in a graph; accepting states get a double circle; start states are explicitly mentioned or indicated otherwise. The transition function is represented by the edges: whenever $d Q_i x$ is a state Q_j , then there is an arrow labelled x from Q_i to Q_j . The input alphabet is implicit in the labels. For automaton M_0 above, the pictorial representation is:



Note that d is a partial function: for example $d B a$ is not defined. We can make d into a total function by introducing a new ‘sink’ state, the result state of all undefined transitions. For example, in the above automaton we can introduce a sink state D with $d D x = D$ for all terminals x , and $d E x = D$ for all states E and terminals x for which $d E x$ is undefined. The sink state and the transitions from/to it are almost always omitted.

The action of a DFA on an input string is described as follows: given a sequence w of input symbols, w can be ‘processed’ symbol by symbol (from left to right) and — depending on the specific input symbol — the DFA (initially in the start state) moves to the state as determined by its state transition function. If no move is possible, the automaton blocks. When the complete input has been processed and the DFA is in one of its accepting states, then we say that w is *accepted by the automaton*.

accept

To illustrate the action of an DFA, we will show that the sentence **bac** is accepted by M_0 . We do so by recording the successive configurations, i.e. the pairs of current state and remaining input values.

$$\begin{aligned}
 & (S, \text{bac}) \\
 \mapsto & \\
 & (A, \text{ac}) \\
 \mapsto & \\
 & (B, \text{c}) \\
 \mapsto & \\
 & (C, \epsilon)
 \end{aligned}$$

Because of the deterministic behaviour of a DFA the definition of acceptance by a DFA is relatively easy. Informally, a sequence $w \in X^*$ is accepted by a DFA (X, Q, d, S, F) , if it is possible, when starting the DFA in S , to end in an accepting state after processing w . This operational description of acceptance is formalised in the predicate *dfa_accept*. The predicate will be derived in a top-down fashion, i.e. we formulate the predicate in terms of (“smaller”) subcomponents and afterwards we give solutions to the subcomponents.

8. Regular Languages

Suppose dfa is a function that reflects the behaviour of the DFA, i.e. a function which given a transition function, a start state and a string, returns the unique state that is reached after processing the string starting from the start state. Then the predicate dfa_accept is defined by:

$$\begin{aligned} dfa_accept &:: X^* \rightarrow (Q \rightarrow X \rightarrow Q, Q, \{Q\}) \rightarrow Bool \\ dfa_accept\ w\ (d, S, F) &= (dfa\ d\ S\ w) \in F \end{aligned}$$

It remains to construct a definition of function dfa that takes a transition function, a start state, and a list of input symbols, and reflects the behaviour of a DFA. The definition of dfa is straightforward

$$\begin{aligned} dfa &:: (Q \rightarrow X \rightarrow Q) \rightarrow Q \rightarrow X^* \rightarrow Q \\ dfa\ d\ q\ \epsilon &= q \\ dfa\ d\ q\ (ax) &= dfa\ d\ (d\ q\ a)\ x \end{aligned}$$

Note that both the type and the definition of dfa match the pattern of the function $foldl$, and it follows that the function dfa is actually identical to $foldl$.

$$dfa\ d\ q = foldl\ d\ q.$$

Definition 8.2 (Acceptance by a DFA). The sequence $w \in X^*$ is accepted by DFA (X, Q, d, S, F) if

$$dfa_accept\ w\ (d, S, F)$$

where

$$\begin{aligned} dfa_accept\ w\ (d, qs, fs) &= dfa\ d\ qs\ w \in fs \\ dfa\ d\ qs &= foldl\ d\ qs \end{aligned}$$

Using the predicate dfa_accept , the language of a DFA is defined as follows.

Definition 8.3 (Language of a DFA). For DFA $M = (X, Q, d, S, F)$, the language of M , $Ldfa(M)$, is defined by

$$Ldfa(M) = \{w \in X^* \mid dfa_accept\ w\ (d, S, F)\}$$

8.1.2. Nondeterministic finite-state automata

This subsection introduces nondeterministic finite-state automata and defines their semantics, i.e. the language of a nondeterministic finite-state automaton.

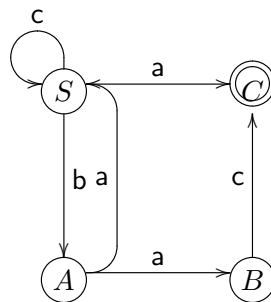
The transition function of a DFA returns a state, which implies that for all terminal symbols x and for all states t there can only be one edge starting in t labelled with x . Sometimes it is convenient to have two or more edges labelled with the same terminal symbol from a state. In these cases one can use a *nondeterministic finite-state automaton*. Nondeterministic finite state automata are defined as follows.

Definition 8.4 (Nondeterministic finite-state automaton, NFA). A nondeterministic finite-state automaton (NFA) is a 5-tuple (X, Q, d, Q_0, F) , where

- X is the input alphabet,
- Q is a finite set of states,
- $d :: Q \rightarrow X \rightarrow \{Q\}$ is the state transition function,
- $Q_0 \subseteq Q$ is the set of start states,
- $F \subseteq Q$ is the set of accepting states.

nondeterministic
finite-state
automaton

An NFA differs from a DFA in that there may be more than one start state and that there may be more than one possible move for each state and input symbol. Here is an example of an NFA:



Note that this NFA is very similar to the DFA in the previous section: the only difference is that there are two outgoing arrows labelled with a from state A . Thus the DFA becomes an NFA.

Formally, this NFA is defined as $M_1 = (X, Q, d, Q_0, F)$ with

$$\begin{aligned} X &= \{a, b, c\} \\ Q &= \{S, A, B, C\} \\ Q_0 &= \{S\} \\ F &= \{C\} \end{aligned}$$

where state transition function d is defined by

$$\begin{aligned} d S a &= \{C\} \\ d S b &= \{A\} \\ d S c &= \{S\} \\ d A a &= \{S, B\} \\ d B c &= \{C\} \end{aligned}$$

Again d is a partial function, which can be made total by adding $d D x = \{ \}$ for all states D and all terminal symbols x for which $d D x$ is undefined.

8. Regular Languages

Since an NFA can make an arbitrary (nondeterministic) choice for one of its possible moves, we have to be careful in defining what it means that a sequence is accepted by an NFA. Informally, sequence $w \in X^*$ is accepted by NFA (X, Q, d, Q_0, F) , if it is possible, when starting the NFA in a state from Q_0 , to end in an accepting state after processing w . This operational description of acceptance is formalised in the predicate *nfa_accept*.

Assume that we have a function, say *nfa*, which reflects the behaviour of the NFA. That is a function which given a transition function, a set of start states and a string, returns all possible states that can be reached after processing the string starting in some start state. Then the predicate *nfa_accept* can be expressed as

$$\begin{aligned} \textit{nfa_accept} &:: X^* \rightarrow (Q \rightarrow X \rightarrow \{Q\}, \{Q\}, \{Q\}) \rightarrow \textit{Bool} \\ \textit{nfa_accept } w (d, Q_0, F) &= \textit{nfa } d Q_0 w \cap F \neq \emptyset \end{aligned}$$

Now it remains to find a function *nfa d qs* of type $X^* \rightarrow \{Q\}$ that reflects the behaviour of the NFA. For lists of length 1 such a function, called *deltas*, is defined by

$$\begin{aligned} \textit{deltas} &:: (Q \rightarrow X \rightarrow \{Q\}) \rightarrow \{Q\} \rightarrow X \rightarrow \{Q\} \\ \textit{deltas } d qs a &= \{r \mid q \in qs, r \in d q a\} \end{aligned}$$

The behaviour of the NFA on X -sequences of arbitrary length follows from this “one step” behaviour:

$$\begin{aligned} \textit{nfa} &:: (Q \rightarrow X \rightarrow \{Q\}) \rightarrow \{Q\} \rightarrow X^* \rightarrow \{Q\} \\ \textit{nfa } d qs \epsilon &= qs \\ \textit{nfa } d qs (ax) &= \textit{nfa } d (\textit{deltas } d qs a) x \end{aligned}$$

Again, it follows that *nfa* can be written as a *foldl*.

$$\textit{nfa } d qs = \textit{foldl } (\textit{deltas } d) qs$$

This concludes the definition of predicate *nfa_accept*. In summary we have derived

Definition 8.5 (Acceptance by an NFA). The sequence $w \in X^*$ is accepted by NFA (X, Q, d, Q_0, F) if

$$\textit{nfa_accept } w (d, Q_0, F)$$

where

$$\begin{aligned} \textit{nfa_accept } w (d, qs, fs) &= \textit{nfa } d qs w \cap fs \neq \emptyset \\ \textit{nfa } d qs &= \textit{foldl } (\textit{deltas } d) qs \\ \textit{deltas } d qs a &= \{r \mid q \in qs, r \in d q a\} \end{aligned}$$

Using the *nfa_accept*-predicate, the language of an NFA is defined by

Definition 8.6 (Language of an NFA). For NFA $M = (X, Q, d, Q_0, F)$, the language of M , $Lnfa(M)$, is defined by

$$Lnfa(M) = \{w \in X^* \mid nfa_accept\ w\ (d, Q_0, F)\}$$

Note that it is computationally expensive to determine whether or not a list is an element of the language of a nondeterministic finite-state automaton. This is due to the fact that all possible transitions have to be tried in order to determine whether or not the automaton can end in an accepting state after reading the input. Determining whether or not a list is an element of the language of a deterministic finite-state automaton can be done in time linear in the length of the input list, so from a computational view, deterministic finite-state automata are preferable. Fortunately, for each nondeterministic finite-state automaton there exists a deterministic finite-state automaton that accepts the same language. We will show how to construct a DFA from an NFA in subsection 8.1.4.

8.1.3. Implementation

This section describes how to implement finite state machines. We start with implementing DFA's. Given a DFA $M = (X, Q, d, S, F)$, we define two datatypes:

```
data StateM   = ... deriving Eq
data SymbolM  = ...
```

where the states of M (the elements of the set Q) are listed as constructors of `StateM`, and the symbols of M (the elements of the set X) are listed as constructors of `SymbolM`. Furthermore, we define three values (one of which a function):

```
start  :: StateM
delta  :: SymbolM -> StateM -> StateM
finals :: [StateM]
```

Note that the first two arguments of `delta` have changed places: this has been done in order to be able to apply 'partial evaluation' later. The extended transition function `dfa` and the accept function `dfaAccept` are now defined by:

```
dfa  :: [SymbolM] -> StateM
dfa  = foldl (flip delta) start

dfaAccept  :: [SymbolM] -> Bool
dfaAccept xs = elem (dfa xs) finals
```

8. Regular Languages

Given a list of symbols $[x_1, x_2, \dots, x_n]$, the computation of `dfa [x1,x2,...,xn]` uses the following intermediate states:

```
start, delta x1 start, delta x2 (delta x1 start),...
```

This list of states is determined uniquely by the input $[x_1, x_2, \dots, x_n]$ and the start state.

Since we want to use the same function names for different automata, we introduce the following class:

```
class Eq a => DFA a b where
  start  :: a
  delta  :: b -> a -> a
  finals :: [a]

  dfa  :: [b] -> a
  dfa = foldl (flip delta) start

  dfaAccept  :: [b] -> Bool
  dfaAccept xs = elem (dfa xs) finals
```

Note that the functions `dfa` and `dfaAccept` are defined once and for all for all instances of the class `DFA`.

As an example, we give the implementation of the example DFA (called `MEX` here) given in the previous subsection.

```
data StateMEX = A | B | C | S deriving Eq
data SymbolMEX = SA | SB | SC
```

So the state `A` is represented by `A`, and the symbol `a` is represented by `SA`, and similar for the other states and symbols. The automaton is made an instance of class `DFA` as follows:

```
instance DFA StateMEX SymbolMEX where
  start = S

  delta x S = case x of SA -> C
                SB -> A
                SC -> S

  delta SA A = B
  delta SC B = C

  finals = [C]
```

We can improve the performance of the automaton (function) `dfa` by means of *partial evaluation*. The main idea of partial evaluation is to replace computations that are performed often at run-time by a single computation that is performed only once at compile-time. A very simple example is the replacement of the expression `if True then f1 else f2` by the expression `f1`. Partial evaluation applies to finite automata in the following way.

partial evaluation

```

dfa [x1,x2,...,xn]
=
foldl (flip delta) start [x1,x2,...,xn]
=
foldl (flip delta) (delta x1 start) [x2,...,xn]
=
case x1 of
  SA -> foldl (flip delta) (delta SA start) [x2,...,xn]
  SB -> foldl (flip delta) (delta SB start) [x2,...,xn]
  SC -> foldl (flip delta) (delta SC start) [x2,...,xn]

```

All these equalities are simple transformation steps for functional programs. Note that the first argument of `foldl` is always `flip delta`, and the second argument is one of the four states S, A, B, or C (the result of `delta`). Since there are only a finite number of states (four, to be precise), we can define a transition function for each state:

```
dfaS, dfaA, dfaB, dfaC :: [Symbol] -> State
```

Each of these functions is a case expression over the possible input symbols.

```

dfaS []      = S
dfaS (x:xs)  = case x of SA -> dfaC xs
                SB -> dfaA xs
                SC -> dfaS xs

dfaA []      = A
dfaA (x:xs)  = case x of SA -> dfaB xs

dfaB []      = B
dfaB (x:xs)  = case x of SC -> dfaC xs

dfaC []      = C

```

With this definition of the finite automaton, the number of steps required for computing the value of `dfaS xs` for some list of symbols `xs` is reduced considerably.

8. Regular Languages

The implementation of NFA's is similar to the implementation of DFA's. The only difference is that the transition and accept functions have to take care of sets (lists) of states now. We will use the following class, in which we use some names that also appear in the class DFA. This is a problem if the two classes appear in the same module.

```
class Eq a => NFA a b where
  start    :: [a]
  delta    :: b -> a -> [a]
  finals   :: [a]

  nfa      :: [b] -> [a]
  nfa      = foldl (flip deltas) start

  deltas   :: b -> [a] -> [a]
  deltas a = union . map (delta a)

  nfaAccept :: [b] -> Bool
  nfaAccept xs = intersect (nfa xs) finals /= []
```

Here, functions `union` and `intersect` are implemented as follows:

```
union :: Eq a => [[a]] -> [a]
union = nub . concat

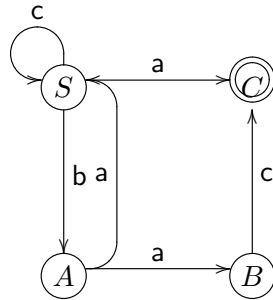
nub :: Eq a => [a] -> [a]
nub = foldr (\x xs -> x:filter (/=x) xs) []

intersect :: Eq a => [a] -> [a] -> [a]
intersect xs ys = intersect' (nub xs)
  where intersect' =
        foldr (\x xs -> if x `elem` ys then x:xs else xs) []
```

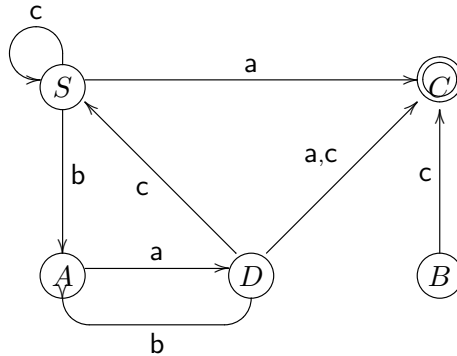
8.1.4. Constructing a DFA from an NFA

Is it possible to express more languages by means of nondeterministic finite-state automata than by deterministic finite-state automata? For each nondeterministic automaton it is possible to give a deterministic finite-state automaton such that both automata accept the same language, so the answer to the above question is no. Before we give the formal proof of this claim, we illustrate the construction of a DFA for an NFA in an example.

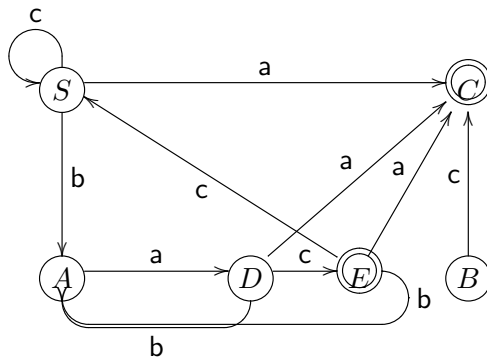
Consider the nondeterministic finite-state automaton corresponding with the example grammar of the previous subsection.



The nondeterminicity of this automaton appears in state A : two outgoing arcs of A are labelled with an a . Suppose we add a new state D , with an arc from A to D labelled a , and we remove the arcs labelled a from A to S and from A to B . Since D is a merge of the states S and B we have to merge the outgoing arcs from S and B into outgoing arcs of D . We obtain the following automaton.



We omit the proof that the language of the latter automaton is equal to the language of the former one. Although there is just one outgoing arc from A labelled with a , this automaton is still nondeterministic: there are two outgoing arcs labelled with c from D . We apply the same procedure as above. Add a new state E and an arc labelled c from D to E , and remove the two outgoing arcs labelled c from D . Since E is a merge of the states C and S we have to merge the outgoing arcs from C and S into outgoing arcs of E . We obtain the following automaton.



8. Regular Languages

Again, we do not prove that the language of this automaton is equal to the language of the previous automaton, provided we add the state E to the set of accepting states, which until now consisted just of state C . State E is added to the set of accepting states because it is the merge of a set of states among which at least one belongs to the set of accepting states. Note that in this automaton for each state, all outgoing arcs are labelled differently, i.e. this automaton is deterministic. The DFA constructed from the NFA above is the 5-tuple (X, Q, d, S, F) with

$$\begin{aligned} X &= \{a, b, c\} \\ Q &= \{S, A, B, C, D, E\} \\ F &= \{C, E\} \end{aligned}$$

where transition function d is defined by

$$\begin{aligned} d S a &= C \\ d S b &= A \\ d S c &= S \\ d A a &= D \\ d B c &= C \\ d D a &= C \\ d D b &= A \\ d D c &= E \\ d E a &= C \\ d E b &= A \\ d E c &= S \end{aligned}$$

This construction is called the ‘subset construction’. In general, the construction works as follows. Suppose $M = (X, Q, d, Q_0, F)$ is a nondeterministic finite-state automaton. Then the finite-state automaton $M' = (X', Q', d', Q'_0, F')$, the components of which are defined below, accepts the same language.

$$\begin{aligned} X' &= X \\ Q' &= \text{subs } Q \end{aligned}$$

where subs returns all subsets of a set. $\text{subs } Q$ is also called the powerset of Q . For example,

$$\begin{aligned} \text{subs } \{X\} &:: \{X\} \rightarrow \{\{X\}\} \\ \text{subs } \{A, B\} &= \{\{\}, \{A\}, \{A, B\}, \{B\}\} \end{aligned}$$

For the other components of M' we define

$$\begin{aligned} d' q a &= \{t \mid t \in d r a, r \in q\} \\ Q'_0 &= \{Q_0\} \\ F' &= \{p \mid p \cap F \neq \emptyset, p \in Q'\} \end{aligned}$$

The proof of the following theorem is given in Section 8.4.

Theorem 8.7 (DFA for NFA). *For every nondeterministic finite-state automaton M there exists a finite-state automaton M' such that*

$$Lnfa(M) = Ldfa(M')$$

Theorem 8.7 enables us to freely switch between NFA's and DFA's. Equipped with this knowledge we continue the exploration of regular languages in the following section. But first we show that the transformation from an NFA to a DFA is an instance of partial evaluation.

8.1.5. Partial evaluation of NFA's

Given a nondeterministic finite state automaton we can obtain a deterministic finite state automaton not just by means of the above construction, but also by means of partial evaluation.

Just as for function `dfa`, we can calculate as follows with function `nfa`.

```
nfa [x1,x2,...,xn]
=
foldl (flip deltas) start [x1,x2,...,xn]
=
foldl (flip deltas) (deltas x1 start) [x2,...,xn]
=
case x1 of
  SA -> foldl (flip deltas) (deltas SA start) [x2,...,xn]
  SB -> foldl (flip deltas) (deltas SB start) [x2,...,xn]
  SC -> foldl (flip deltas) (deltas SC start) [x2,...,xn]
```

Note that the first argument of `foldl` is always `flip deltas`, and the second argument is one of the six sets of states [S], [A], [B], [C], [B,S], [C,S] (the possible results of `deltas`). Since there are only a finite number of results of `deltas` (six, to be precise), we can define a transition function for each state:

```
nfaS, nfaA, nfaB, nfaC, nfaBS, nfaCS :: [Symbol] -> [State]
```

For example,

```
nfaA [] = A
nfaA (x:xs) = case x of
  SA -> nfaBS xs
  _ -> error "no transition possible"
```

8. Regular Languages

Each of these functions is a case expression over the possible input symbols. By partially evaluating the function `nfa` we have obtained a function that is the implementation of the deterministic finite state automaton corresponding to the nondeterministic finite state automaton.

8.2. Regular grammars

This section defines regular grammars, a special kind of context-free grammars. Subsection 8.2.1 gives the correspondence between nondeterministic finite-state automata and regular grammars.

regular grammar

Definition 8.8 (Regular Grammar). A *regular grammar* G is a context free grammar (T, N, R, S) in which all production rules in R are of one of the following two forms:

$$\begin{aligned} A &\rightarrow xB \\ A &\rightarrow x \end{aligned}$$

with $x \in T^*$ and $A, B \in N$. So in every rule there is at most one nonterminal, and if there is a nonterminal present, it occurs at the end.

The regular grammars as defined here are sometimes called right-regular grammars. There is a symmetric definition for left-regular grammars.

regular language

Definition 8.9 (Regular Language). A *regular language* is a language that is generated by a regular grammar.

From the definition it is clear that each regular language is context-free. The question is now: Is each context-free language regular? The answer is: No. There are context-free languages that are not regular; an example of such a language is $\{a^n b^n \mid n \in \mathbb{N}\}$. To understand this, you have to know how to prove that a language is not regular. Because of its subtlety, we postpone this kind of proofs until Chapter 9. Here it suffices to know that regular languages form a proper subset of context-free languages and that we will profit from their speciality in the recognition process.

A first similarity between regular languages and context-free languages is that both are closed under union, concatenation and Kleene-star.

Theorem 8.10. *Let L and M be regular languages, then*

$$\begin{aligned} L \cup M &\text{ is regular} \\ LM &\text{ is regular} \\ L^* &\text{ is regular} \end{aligned}$$

Proof. Let $G_L = (T, N_L, R_L, S_L)$ and $G_M = (T, N_M, R_M, S_M)$ be regular grammars for L and M respectively, then

- for regular grammars, the well-known union construction for context-free grammars is a regular grammar again;
- we obtain a regular grammar for LM if we replace, in G_L , each production of the form $T \rightarrow x$ and $T \rightarrow \epsilon$ by $T \rightarrow xS_M$ and $T \rightarrow S_M$, respectively;
- since $L^* = \{\epsilon\} \cup LL^*$, it follows from the above that there exists a regular grammar for L^* .

□

In addition to these closure properties, regular languages are closed under intersection and complement too. See the exercises. This is remarkable because context-free languages are not closed under these operations. Recall the language $L = L_1 \cap L_2$ where $L_1 = \{a^n b^n c^m \mid n, m \in \mathbb{N}\}$ and $L_2 = \{a^n b^m c^m \mid n, m \in \mathbb{N}\}$.

As for context-free languages, there may exist more than one regular grammar for a given regular language and these regular grammars may be transformed into each other.

We conclude this section with a grammar transformation:

Theorem 8.11. *For each regular grammar G there exists a regular grammar G' with start-symbol S' such that*

$$L(G) = L(G')$$

and such that G' has no productions of the form $U \rightarrow V$ and $W \rightarrow \epsilon$, with V a single nonterminal and $W \neq S'$.

In other words: every regular grammar can be transformed to a form where every production has a nonempty terminal string in its right hand side (with a possible exception for $S \rightarrow \epsilon$).

The proof of this transformation is omitted, we only briefly describe the construction of such a regular grammar, and illustrate the construction with an example.

Given a regular grammar G , a regular grammar with the same language but without productions of the form $U \rightarrow V$ and $W \rightarrow \epsilon$ for all U, V , and all $W \neq S$ is obtained as follows. First, consider all pairs Y, Z of nonterminals of G such that $Y \xrightarrow{*} Z$. Add productions $Y \rightarrow z$ to the grammar, with $Z \rightarrow z$ a production of the original grammar, and z not a single nonterminal. Remove all productions $U \rightarrow V$ from G . Finally, remove all productions of the form $W \rightarrow \epsilon$ for $W \neq S$, and for each production $U \rightarrow xW$ add the production $U \rightarrow x$. The following example illustrates this construction.

8. Regular Languages

Consider the following regular grammar G .

$$\begin{aligned} S &\rightarrow aA \\ S &\rightarrow bB \\ S &\rightarrow A \\ S &\rightarrow C \\ A &\rightarrow bB \\ A &\rightarrow S \\ A &\rightarrow \epsilon \\ B &\rightarrow bB \\ B &\rightarrow \epsilon \\ C &\rightarrow c \end{aligned}$$

The grammar G' of the desired form is constructed in 3 steps.

Step 1

Let G' equal G .

Step 2

Consider all pairs of nonterminals Y and Z . If $Y \xRightarrow{*} Z$, add the productions $Y \rightarrow z$ to G' , with $Z \rightarrow z$ a production of the original grammar, and z not a single nonterminal. Furthermore, remove all productions of the form $U \rightarrow V$ from G' . In the example we remove the productions $S \rightarrow A$, $S \rightarrow C$, $A \rightarrow S$, and we add the productions $S \rightarrow bB$ and $S \rightarrow \epsilon$ since $S \xRightarrow{*} A$, and the production $S \rightarrow c$ since $S \xRightarrow{*} C$, and the productions $A \rightarrow aA$ and $A \rightarrow bB$ since $A \xRightarrow{*} S$, and the production $A \rightarrow c$ since $A \xRightarrow{*} C$. We obtain the grammar with the following productions.

$$\begin{aligned} S &\rightarrow aA \\ S &\rightarrow bB \\ S &\rightarrow c \\ S &\rightarrow \epsilon \\ A &\rightarrow bB \\ A &\rightarrow aA \\ A &\rightarrow c \\ A &\rightarrow \epsilon \\ B &\rightarrow bB \\ B &\rightarrow \epsilon \\ C &\rightarrow c \end{aligned}$$

This grammar generates the same language as G , and has no productions of the form $U \rightarrow V$. It remains to remove productions of the form $W \rightarrow \epsilon$ for $W \neq S$.

Step 3

Remove all productions of the form $W \rightarrow \epsilon$ for $W \neq S$, and for each production $U \rightarrow xW$ add the production $U \rightarrow x$. Applying this transformation to the above grammar gives the following grammar.

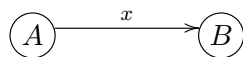
$$\begin{aligned}
 S &\rightarrow aA \\
 S &\rightarrow bB \\
 S &\rightarrow a \\
 S &\rightarrow b \\
 S &\rightarrow c \\
 S &\rightarrow \epsilon \\
 A &\rightarrow bB \\
 A &\rightarrow aA \\
 A &\rightarrow a \\
 A &\rightarrow b \\
 A &\rightarrow c \\
 B &\rightarrow bB \\
 B &\rightarrow b \\
 C &\rightarrow c
 \end{aligned}$$

Each production in this grammar is of one of the desired forms: $U \rightarrow x$ or $U \rightarrow xV$, and the language of the grammar G' we thus obtain is equal to the language of grammar G .

8.2.1. Equivalence of Regular grammars and Finite automata

In the previous section, we introduced finite-state automata. Here we show that regular grammars and nondeterministic finite-state automata are two sides of one coin.

We will prove the equivalence using theorem 8.7. The equivalence consists of two parts, formulated in the theorems 8.12 and 8.13 below. The basis for both theorems is the direct correspondence between a production $A \rightarrow xB$ and a transition



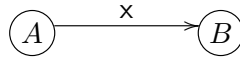
Theorem 8.12 (Regular grammar for NFA). *For each NFA M there exists a regular grammar G such that*

$$L_{nfa}(M) = L(G)$$

8. Regular Languages

Proof. We will just sketch the construction, the formal proof can be found in the literature. Let (X, Q, d, S, F) be a DFA for NFA M . Construct the grammar $G = (X, Q, R, S)$ where

- the terminals of the grammar are the input alphabet of the automaton;
- the nonterminals of the grammar are the states of the automaton;
- the start state of the grammar is the start state of the automaton;
- the productions of the grammar correspond to the automaton transitions:
a rule $A \rightarrow xB$ for each transition



a rule $A \rightarrow \epsilon$ for each terminal state A .

In formulae:

$$R = \{A \rightarrow xB \mid A, B \in Q, x \in X, d A x = B\} \\ \cup \{A \rightarrow \epsilon \mid A \in F\}$$

□

Theorem 8.13 (NFA for regular grammar). *For each regular grammar G there exists a nondeterministic finite-state automaton M such that*

$$L(G) = L_{nfa}(M)$$

Proof. Again, we will just sketch the construction, the formal proof can be found in the literature. The construction consists of two steps: first we give a direct translation of a regular grammar to an automaton and then we transform the automaton into a suitable shape.

From a grammar to an automaton. Let $G = (T, N, R, S)$ be a regular grammar without productions of the form $U \rightarrow V$ and $W \rightarrow \epsilon$ for $W \neq S$. Construct NFA $M = (X, Q, d, \{S\}, F)$ where

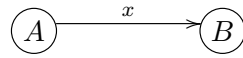
- The input alphabet of the automaton are the nonempty terminal *strings* (!) that occur in the rules of the grammar:

$$X = \{x \in T^+ \mid A, B \in N, A \rightarrow xB \in R\} \\ \cup \{x \in T^+ \mid A \in N, A \rightarrow x \in R\}$$

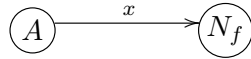
- The states of the automaton are the nonterminals of the grammar extended with a new state N_f .

$$Q = N \cup \{N_f\}$$

- The transitions of the automaton correspond to the grammar productions: for each rule $A \rightarrow xB$ we get a transition



for each rule $A \rightarrow x$ with nonempty x , we get a transition



In formulae: For all $A \in N$ and $x \in X$:

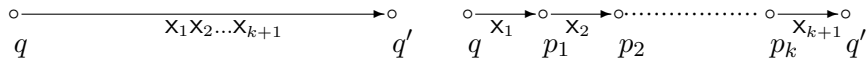
$$\begin{aligned} \delta A x &= \{B \mid B \in N, A \rightarrow xB \in R\} \\ &\cup \{N_f \mid A \rightarrow x \in R\} \end{aligned}$$

- The final states of the automaton are N_f and possibly S , if $S \rightarrow \epsilon$ is a grammar production.

$$F = \{N_f\} \cup \{S \mid S \rightarrow \epsilon \in R\}$$

$L_{nfa}(M) = L(G)$, because of the direct correspondence between derivation steps in G and transitions in M .

Transforming the automaton to a suitable shape. There is a minor flaw in the automaton given above: the grammar and the automaton have different alphabets. This shortcoming will be remedied by an automaton transformation which yields an equivalent automaton with transitions labelled by elements of T (instead of T^*). The transformation is relatively easy and is depicted in the diagram below. In order to eliminate transition $\delta q x = q'$ where $x = x_1x_2 \dots x_{k+1}$ with $k > 0$ and $x_i \in T$ for all i , add new (nonfinal) states p_1, \dots, p_k to the existing ones and new transitions $\delta q x_1 = p_1, \delta p_1 x_2 = p_2, \dots, \delta p_k x_{k+1} = q'$.



It is intuitively clear that the resulting automaton is equivalent to the original one. Carry out this transformation for each M -transition $\delta q x = q'$ with $|x| > 1$ in order to get an automaton for G with the same input alphabet T . \square

8.3. Regular expressions

Regular expressions are a classical and convenient way to describe, for example, the structure of terminal words. This section defines regular expressions, defines the language of a regular expression, and shows that regular expressions and regular grammars are equally expressive formalisms. We do not discuss implementations of (datatypes and functions for matching) regular expressions; implementations can be found in the literature [8, 6].

Definition 8.14 (RE_T , regular expressions over alphabet T). The set RE_T of regular expressions over alphabet T is inductively defined as follows: for regular expressions R, S

regular expression

$$\begin{aligned}
 \emptyset &\in RE_T \\
 \epsilon &\in RE_T \\
 \mathbf{a} &\in RE_T \\
 R + S &\in RE_T \\
 RS &\in RE_T \\
 R^* &\in RE_T \\
 (R) &\in RE_T
 \end{aligned}$$

where $\mathbf{a} \in T$. The operator $+$ is associative, commutative, and idempotent; the concatenation operator, written as juxtaposition (so x concatenated with y is denoted by xy), is associative, and ϵ is the unit of it. In formulae this reads, for all regular expressions R, S , and V ,

$$\begin{aligned}
 R + (S + U) &= (R + S) + U \\
 R + S &= S + R \\
 R + R &= R
 \end{aligned}$$

$$\begin{aligned}
 R(SU) &= (RS)U \\
 R\epsilon &= R \quad (= \epsilon R)
 \end{aligned}$$

Furthermore, the star operator, $*$, binds stronger than concatenation, and concatenation binds stronger than $+$. Examples of regular expressions are:

$$\begin{aligned}
 (\mathbf{bc})^* + \emptyset \\
 \epsilon + \mathbf{b}(\epsilon^*)
 \end{aligned}$$

The language (i.e. the “semantics”) of a regular expression over T is a set of T -sequences compositionally defined on the structure of regular expressions. As follows.

Definition 8.15 (Language of a regular expression). Function $Lre :: RE_T \rightarrow \{T^*\}$ returns the language of a regular expression. It is defined inductively by:

$$\begin{aligned}
 Lre(\emptyset) &= \emptyset \\
 Lre(\epsilon) &= \{\epsilon\} \\
 Lre(\mathbf{b}) &= \{\mathbf{b}\} \\
 Lre(x + y) &= Lre(x) \cup Lre(y) \\
 Lre(xy) &= Lre(x) Lre(y) \\
 Lre(x^*) &= (Lre(x))^*
 \end{aligned}$$

Since \cup is associative, commutative, and idempotent, set concatenation is associative with $\{\epsilon\}$ as its unit, and function Lre is well defined. Note that the language $Lre\mathbf{b}^*$ is the set consisting of zero, one or more concatenations of \mathbf{b} , i.e., $Lre(\mathbf{b}^*) = (\{\mathbf{b}\})^*$. As an example of a language of a regular expression, we compute the language of the regular expression $(\epsilon + \mathbf{bc})\mathbf{d}$.

$$\begin{aligned}
 &Lre((\epsilon + \mathbf{bc})\mathbf{d}) \\
 = & \\
 &(Lre(\epsilon + \mathbf{bc}))(Lre(\mathbf{d})) \\
 = & \\
 &(Lre(\epsilon) \cup Lre(\mathbf{bc}))\{\mathbf{d}\} \\
 = & \\
 &(\{\epsilon\} \cup (Lre(\mathbf{b}))(Lre(\mathbf{c})))\{\mathbf{d}\} \\
 = & \\
 &\{\epsilon, \mathbf{bc}\}\{\mathbf{d}\} \\
 = & \\
 &\{\mathbf{d}, \mathbf{bcd}\}
 \end{aligned}$$

Regular expressions are used to describe the *tokens* of a language. For example, the list

if p then e1 else e2

contains six tokens, three of which are identifiers. An *identifier* is an element in the language of the regular expression

$$letter(letter + digit)^*$$

where

$$\begin{aligned}
 letter &= \mathbf{a} + \mathbf{b} + \dots + \mathbf{z} + \\
 &\quad \mathbf{A} + \mathbf{B} + \dots + \mathbf{Z} \\
 digit &= \mathbf{0} + \mathbf{1} + \dots + \mathbf{9}
 \end{aligned}$$

8. Regular Languages

see subsection 2.3.1.

In the beginning of this section we claimed that regular expressions and regular grammars are equivalent formalisms. We will prove this claim later, but first we illustrate the construction of a regular grammar out of a regular expressions in an example. Consider the following regular expression.

$$R = a^* + \epsilon + (a + b)^*$$

We aim at a regular grammar G such that $Lre(R) = L(G)$ and again we take a top-down approach.

Suppose that nonterminal A generates the language $Lre(a^*)$, nonterminal B generates the language $Lre(\epsilon)$, and nonterminal C generates the language $Lre((a + b)^*)$. Suppose furthermore that the productions for A , B , and C satisfy the conditions imposed upon regular grammars. Then we obtain a regular grammar G with $L(G) = Lre(R)$ by defining

$$\begin{aligned} S &\rightarrow A \\ S &\rightarrow B \\ S &\rightarrow C \end{aligned}$$

where S is the start-symbol of G . It remains to construct productions for nonterminals A , B , and C .

- The nonterminal A with productions

$$\begin{aligned} A &\rightarrow aA \\ A &\rightarrow \epsilon \end{aligned}$$

generates the language $Lre(a^*)$.

- Since $Lre(\epsilon) = \{\epsilon\}$, the nonterminal B with production

$$B \rightarrow \epsilon$$

generates the language $\{\epsilon\}$.

- Nonterminal C with productions

$$\begin{aligned} C &\rightarrow aC \\ C &\rightarrow bC \\ C &\rightarrow \epsilon \end{aligned}$$

generates the language $Lre((a + b)^*)$.

For a specific example it is not difficult to construct a regular grammar for a regular expression. We now give the general result.

Theorem 8.16 (Regular Grammar for Regular Expression). *For each regular expression R there exists a regular grammar G such that*

$$Lre(R) = L(G)$$

The proof of this theorem is given in Section 8.4.

To obtain a regular expression that generates the same language as a given regular grammar we go via an automaton. Given a regular grammar G , we can use the theorems from the previous sections to obtain a DFA D such that

$$L(G) = Ldfa(D)$$

So if we can obtain a regular expression for a DFA D , we have found a regular expression for a regular grammar. To obtain a regular expression for a DFA D , we interpret each state of D as a regular expression defined as the sum of the concatenation of outgoing terminal symbols with the resulting state. For our example DFA we obtain:

$$\begin{aligned} S &= aC + bA + cS \\ A &= aB \\ B &= cC \\ C &= \epsilon \end{aligned}$$

It is easy to merge these four regular expressions into a single regular expression, partially because this is a simple example. Merging the regular expressions obtained from a DFA that may loop is more complicated, as we will briefly explain in the proof of the following theorem. In general, we have:

Theorem 8.17 (Regular Expression for Regular Grammar). *For each regular grammar G there exists a regular expression R such that*

$$L(G) = Lre(R)$$

The proof of this theorem is given in Section 8.4.

8.4. Proofs

This section contains the proofs of some of the theorems given in this chapter.

8.4.1. Proof of Theorem 8.7

Suppose $M = (X, Q, d, Q_0, F)$ is a nondeterministic finite-state automaton. Define the finite-state automaton $M' = (X', Q', d', Q'_0, F')$ as follows.

$$\begin{aligned} X' &= X \\ Q' &= \text{subs } Q \end{aligned}$$

where *subs* returns the powerset of a set. For example,

$$\text{subs } \{A, B\} = \{\{\}, \{A\}, \{A, B\}, \{B\}\}$$

For the other components of M' we define

$$\begin{aligned} d' q a &= \{t \mid t \in d r a, r \in q\} \\ Q'_0 &= Q_0 \\ F' &= \{p \mid p \cap F \neq \emptyset, p \in Q'\} \end{aligned}$$

We have

$$\begin{aligned} &Lnfa(M) \\ = &\text{definition of } Lnfa \\ &\{w \mid w \in X^*, nfa_accept w (d, Q_0, F)\} \\ = &\text{definition of } nfa_accept \\ &\{w \mid w \in X^*, (nfa d Q_0 w \cap F) \neq \emptyset\} \\ = &\text{definition of } F' \\ &\{w \mid w \in X^*, nfa d Q_0 w \in F'\} \\ = &\mathbf{assume} \ nfa d Q_0 w = dfa d' Q'_0 w \\ &\{w \mid w \in X^*, dfa d' Q'_0 w \in F'\} \\ = &\text{definition of } dfa_accept \\ &\{w \mid w \in X^*, dfa_accept w (d', Q'_0, F')\} \\ = &\text{definition of } Ldfa \\ &Ldfa(M') \end{aligned}$$

It follows that $Lnfa(M) = Ldfa(M')$ provided

$$nfa d Q_0 w = dfa d' Q'_0 w \tag{8.1}$$

We prove this equation as follows.

$$\begin{aligned}
& nfa\ d\ Q_0\ w \\
= & \text{definition of } nfa \\
& foldl\ (deltas\ d)\ Q_0\ w \\
= & Q_0 = Q'_0 ; \mathbf{assume}\ deltas\ d = d' \\
& foldl\ d'\ Q'_0\ w \\
= & \text{definition of } dfa \\
& dfa\ d'\ Q'_0\ w
\end{aligned}$$

So if

$$deltas\ d = d'$$

equation (8.1) holds. This equality follows from the following calculation.

$$\begin{aligned}
& d'\ q\ \mathbf{a} \\
= & \text{definition of } d' \\
& \{t \mid r \in q, t \in d\ r\ \mathbf{a}\} \\
= & \text{definition of } deltas \\
& deltas\ d\ q\ \mathbf{a}
\end{aligned}$$

8.4.2. Proof of Theorem 8.16

The proof of this theorem is by induction to the structure of regular expressions. For the three base cases we argue as follows.

The regular grammar without productions generates the language $Lre(\emptyset)$. The regular grammar with the production $S \rightarrow \epsilon$ generates the language $Lre(\epsilon)$. The regular grammar with production $S \rightarrow \mathbf{b}$ generates the language $Lre(\mathbf{b})$.

For the other three cases the induction hypothesis is that there exist a regular grammar with start-symbol S_1 that generates the language $Lre(x)$, and a regular grammar with start-symbol S_2 that generates the language $Lre(y)$.

We obtain a regular grammar with start-symbol S that generates the language $Lre(x + y)$ by defining

$$\begin{aligned}
S & \rightarrow S_1 \\
S & \rightarrow S_2
\end{aligned}$$

We obtain a regular grammar with start-symbol S that generates the language $Lre(xy)$ by replacing, in the regular grammar that generates the language $Lre(x)$,

8. Regular Languages

each production of the form $T \rightarrow a$ and $T \rightarrow \epsilon$ by $T \rightarrow aS_2$ and $T \rightarrow S_2$, respectively.

We obtain a regular grammar with start-symbol S that generates the language $Lre(x^*)$ by replacing, in the regular grammar that generates the language $Lre(x)$, each production of the form $T \rightarrow a$ and $T \rightarrow \epsilon$ by $T \rightarrow aS$ and $T \rightarrow S$, and by adding the productions $S \rightarrow S_1$ and $S \rightarrow \epsilon$, where S_1 is the start-symbol of the regular grammar that generates the language $Lre(x)$.

8.4.3. Proof of Theorem 8.17

In sections 8.1.4 and 8.2.1 we have shown that there exists a DFA $D = (X, Q, d, S, F)$ such that

$$L(G) = Ldfa(D)$$

So, if we can show that there exists a regular expression R such that

$$Ldfa(D) = Lre(R)$$

then the theorem follows.

Let $D = (X, Q, d, S, F)$ be a DFA such that $L(G) = Ldfa(D)$. We define a regular expression R such that

$$Lre(R) = Ldfa(D)$$

For each state $q \in Q$ we define a regular expression \bar{q} , and we let R be \bar{S} . We obtain the definition of \bar{q} by combining all pairs c and C such that $d q c = C$.

$$\begin{aligned} \bar{q} &= \text{if } q \notin F \\ &\quad \text{then foldl } (+) \emptyset [cC \mid d q c = C] \\ &\quad \text{else } \epsilon + \text{foldl } (+) \emptyset [cC \mid d q c = C] \end{aligned}$$

This gives a set of possibly mutually recursive equations, which we have to solve. In solving these equations we use the fact that concatenation distributes over the sum operator:

$$z(x + y) = zx + zy$$

and that recursion can be removed by means of the star operator $*$:

$$A = xA + z \text{ (where } A \notin z) \equiv A = x^*z$$

The algorithm for solving such a set of equations is omitted.

We prove $Ldfa(D) = Lre(\bar{S})$.

$$\begin{aligned}
& Ldfa(D) \\
= & \text{definition of } Ldfa \\
& \{w \mid w \in X^*, dfa_accept\ w\ (d, S, F)\} \\
= & \text{definition of } dfa_accept \\
& \{w \mid w \in X^*, (dfa\ d\ S\ w) \in F\} \\
= & \text{definition of } dfa \\
& \{w \mid w \in X^*, (foldl\ d\ S\ w) \in F\} \\
= & \text{assumption} \\
& \{w \mid w \in Lre(\bar{S})\} \\
= & \text{equality for set-comprehensions} \\
& Lre(\bar{S})
\end{aligned}$$

It remains to prove the assumption in the above calculation: for $w \in X^*$,

$$(foldl\ d\ S\ w) \in F \equiv w \in Lre(\bar{S})$$

We prove a generalisation of this equation, namely, for arbitrary q ,

$$(foldl\ d\ q\ w) \in F \equiv w \in Lre(\bar{q})$$

This equation is proved by induction to the length of w . For the base case $w = \epsilon$ we calculate as follows.

$$\begin{aligned}
& (foldl\ d\ q\ \epsilon) \in F \\
\equiv & \text{definition of } foldl \\
& q \in F \\
\equiv & \text{definition of } \bar{q}, E \text{ abbreviates the fold expression} \\
& \bar{q} = \epsilon + E \\
\equiv & \text{definition of } Lre, \text{ definition of } \bar{q} \\
& \epsilon \in Lre(\bar{q})
\end{aligned}$$

The induction hypothesis is that for all lists w with $|w| \leq n$ we have $(foldl\ d\ q\ w) \in F \equiv w \in Lre(\bar{q})$. Suppose $\mathbf{a}x$ is a list of length $n+1$.

$$\begin{aligned}
& (foldl\ d\ q\ (\mathbf{a}x)) \in F \\
\equiv & \text{definition of } foldl \\
& (foldl\ d\ (d\ q\ \mathbf{a})\ x) \in F \\
\equiv & \text{induction hypothesis} \\
& x \in Lre(d\ \bar{q}\ \mathbf{a}) \\
\equiv & \text{definition of } \bar{q}, D \text{ is deterministic} \\
& \mathbf{a}x \in Lre(\bar{q})
\end{aligned}$$

Summary

This chapter discusses methods for recognising sentences from regular languages, and introduces several concepts related to describing and recognising regular languages. Regular languages are used for describing simple languages like ‘the language of identifiers’ and ‘the language of keywords’ and regular expressions are convenient for the description of regular languages. The straightforward translation of a regular expression into a recogniser for the language of that regular expression results in a recogniser that is often very inefficient. By means of (non)deterministic finite-state automata we construct a recogniser that requires time linear in the length of the input list for recognising an input list.

8.5. Exercises

Exercise 8.1. Given a regular grammar G for language L , construct a regular grammar for L^* .

Exercise 8.2. Transform the grammar with the following productions to a grammar without productions of the form $U \rightarrow V$ and $W \rightarrow \epsilon$ with $W \neq S$.

$$\begin{aligned} S &\rightarrow aA \\ S &\rightarrow A \\ A &\rightarrow aS \\ A &\rightarrow B \\ B &\rightarrow C \\ B &\rightarrow \epsilon \\ C &\rightarrow cC \\ C &\rightarrow a \end{aligned}$$

Exercise 8.3. Suppose that the state transition function d in the definition of a nondeterministic finite-state automaton has the following type

$$d \quad :: \quad \{Q\} \rightarrow X \rightarrow \{Q\}$$

Function d takes a set of states V and an element a , and returns the set of states that are reachable from V with an arc labelled a . Define a function $ndfsa$ of type

$$(\{Q\} \rightarrow X \rightarrow \{Q\}) \rightarrow \{Q\} \rightarrow X^* \rightarrow \{Q\}$$

which given a function d , a set of start states, and an input list, returns the set of states in which the nondeterministic finite-state automaton can end after reading the input list.

Exercise 8.4. Prove the converse of Theorem 8.7: show that for every deterministic finite-state automaton M there exists a nondeterministic finite-state automaton M' such that

$$Lda(M) = Lna(M')$$

Exercise 8.5. Regular languages are closed under complementation. Prove this claim. Hint: construct a finite automaton for \bar{L} out of an automaton for regular language L .

Exercise 8.6. Regular languages are closed under intersection.

1. Prove this claim using the result from the previous exercise.
2. A direct proof from this claim is the following:

Let $M_1 = (X, Q_1, d_1, S_1, F_1)$ and $M_2 = (X, Q_2, d_2, S_2, F_2)$ be DFA's for the regular languages L_1 and L_2 respectively. Define the (product) automaton $M = (X, Q_1 \times Q_2, d, (S_1, S_2), F_1 \times F_2)$ by $d(q_1, q_2) x = (d_1 q_1 x, d_2 q_2 x)$. Now prove that $Ldfa(M) = Ldfa(M_1) \cap Ldfa(M_2)$.

Exercise 8.7. Define nondeterministic finite-state automata that accept languages equal to the languages of the following regular grammars.

$$1. \begin{cases} S \rightarrow (A \\ S \rightarrow \epsilon \\ S \rightarrow)A \\ A \rightarrow) \\ A \rightarrow (\end{cases}$$

$$2. \begin{cases} S \rightarrow 0A \\ S \rightarrow 0B \\ S \rightarrow 1A \\ A \rightarrow 1 \\ A \rightarrow 0 \\ B \rightarrow 0 \\ B \rightarrow \epsilon \end{cases}$$

Exercise 8.8. Describe the language of the following regular expressions.

1. $\epsilon + b(\epsilon^*)$
2. $(bc)^* + \emptyset$
3. $a(b^*) + c^*$

Exercise 8.9. Prove that for arbitrary regular expressions R , S , and T the following equivalences hold.

$$\begin{aligned} Lre(R(S + T)) &= Lre(RS + RT) \\ Lre((R + S)T) &= Lre(RT + ST) \end{aligned}$$

Exercise 8.10. Give regular expressions S and R such that

$$\begin{aligned} Lre(RS) &= Lre(SR) \\ Lre(RS) &\neq Lre(SR) \end{aligned}$$

Exercise 8.11. Give regular expressions V and W , with $Lre(V) \neq Lre(W)$, such that for all regular expressions R and S with $S \neq \emptyset$

$$Lre(R(S + V)) = Lre(R(S + W))$$

V and W may be expressed in terms of R and S .

Exercise 8.12. Give a regular expression for the language that consists of all lists of zeros and ones such that the segment 01 occurs nowhere in a list. Examples of sentences of this language are 1110, and 000.

8. Regular Languages

Exercise 8.13. Give regular grammars that generate the language of the following regular expressions.

1. $((a + bb)^* + c)^*$
2. $a^* + b^* + ab$

Exercise 8.14. Give regular expressions of which the language equals the language of the following regular grammars.

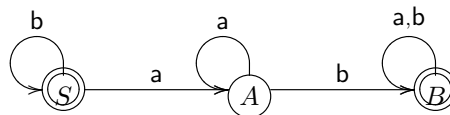
1.
$$\begin{cases} S \rightarrow bA \\ S \rightarrow aC \\ S \rightarrow \epsilon \\ A \rightarrow bA \\ A \rightarrow \epsilon \\ B \rightarrow aC \\ B \rightarrow bB \\ B \rightarrow \epsilon \\ C \rightarrow bB \\ C \rightarrow b \end{cases}$$
2.
$$\begin{cases} S \rightarrow 0S \\ S \rightarrow 1T \\ S \rightarrow \epsilon \\ T \rightarrow 0T \\ T \rightarrow 1S \end{cases}$$

Exercise 8.15. Construct for each of the following regular expressions a nondeterministic finite-state automaton that accepts the sentences of the language. Transform the nondeterministic finite-state automata into deterministic finite-state automata.

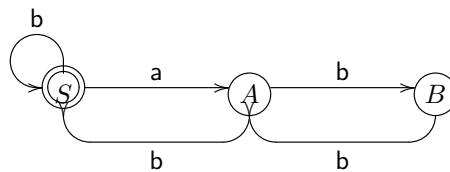
1. $a^* + b^* + (ab)$
2. $(1 + (12) + 0)^*(30)^*$

Exercise 8.16. Define regular expressions for the languages of the following deterministic finite-state automata.

1. Start state is S .



2. Start state is S .



9. Pumping Lemmas: the expressive power of languages

Introduction

In these lecture notes we have presented several ways to show that a language is regular or context-free, but until now we did not give any means to show the non-regularity or noncontext-freeness of a language. In this chapter we fill this gap by introducing the so-called *Pumping Lemmas*. For example, the pumping lemma for regular languages says

IF language L is regular,

THEN it has the following property P : each sufficiently long sentence $w \in L$ has a substring that can be repeated any number of times, every time yielding another word of L

In applications, pumping lemmas are used in the contrapositive way. In the regular case this means that one may conclude that L is *not* regular, if P does *not* hold. Although the ideas behind pumping lemmas are very simple, a precise formulation is not. As a consequence, it takes some effort to get familiar with applying pumping lemmas. Regular grammars and context-free grammars are part of the Chomsky hierarchy, which consists of four different kinds of grammars and their corresponding languages. Pumping lemmas are used to show that the expressive power of the different elements of the Chomsky hierarchy is different.

Goals

After you have studied this chapter you will be able to

- prove that a language is not regular;
- prove that a language is not context-free;
- identify languages and grammars as regular, context-free or none of these;
- give examples of languages that are not regular, and/or not context-free;
- explain the Chomsky hierarchy.

9.1. The Chomsky hierarchy

In the preceding chapters we have seen context-free grammars and regular grammars. You may now wonder: is it possible to express any language with these grammars? And: is it possible to obtain any context-free language from a regular grammar? The answer to these questions is no. The Chomsky hierarchy explains why the answer is no. The Chomsky hierarchy consists of four elements, each of which is explained below.

9.1.1. Type-0 grammars

The most powerful grammars are the type-0 grammars, in which a production has the form $\phi \rightarrow \psi$, where $\phi \in V^+$, $\psi \in V^*$, where V is the set of symbols of the grammar. So the left-hand side of a production may consist of a list of nonterminal and terminal symbols, instead of a single nonterminal as in context-free grammars. Type-0 grammars have the same expressive power as Turing machines, and the languages described by these grammars are the recursive enumerable languages. This expressive power comes at a cost though: it is very difficult to parse sentences from type-0 grammars.

9.1.2. Type-1 grammars

We can slightly restrict the form of the productions to obtain type-1 grammars. In a type-1, or context-sensitive grammar, each production has the form $\phi A \psi \rightarrow \phi \delta \psi$, where $\phi, \psi \in V^*$, $\delta \in V^+$. So a production describes how to rewrite a nonterminal A , in the context of lists of symbols ϕ and ψ . A language generated by a context-sensitive grammar is called a context-sensitive language. Although context-sensitive grammars are less expressive than type-0 grammars, parsing is still very difficult for context-sensitive grammars.

9.1.3. Type-2 grammars

The type-2 grammars are the context-free grammars which you have seen a lot in the preceding chapters. As the name says, in a context-free grammar you can rewrite nonterminals without looking at the context in which they appear. Actually, it is *impossible* to look at the context when rewriting symbols. Context-free grammars are less expressive than context-sensitive grammars. This statement can be proved using the pumping lemma for context-free languages.

However, it is much easier to parse sentences from context-free languages. In fact, a sentence of length n can be parsed in time at most $O(n^3)$ (or even a bit less than this) for any sentence of a context-free language. And if we put some more restrictions

on context-free grammars (for example $LL(1)$), we obtain linear-time algorithms for parsing sentences of such grammars.

9.1.4. Type-3 grammars

The type-3 grammars in the Chomsky hierarchy are the regular grammars. Any sentence from a regular language can be processed by means of a finite-state automaton, which takes linear time and constant space in the size of its input. The set of regular languages is strictly smaller than the set of context-free languages, a fact we will prove below by means of the pumping lemma for regular languages.

9.2. The pumping lemma for regular languages

In this section we give the pumping lemma for regular languages. The lemma gives a property that is satisfied by all regular languages. The property is a statement of the form: in sentences longer than a certain length a substring can be identified that can be duplicated while retaining a sentence. The idea behind this property is simple: regular languages are accepted by finite automata. Given a DFA for a regular language, a sentence of the language describes a path from the start state to some finite state. When the length of such a sentence exceeds the number of states, then at least one state is visited twice; consequently the path contains a cycle that can be repeated as often as desired. The proof of the following lemma is given in Section 9.4.

Theorem 9.1 (Regular Pumping Lemma). *Let L be a regular language. Then*

there exists $n \in \mathbb{N}$:

for all x, y, z : $xyz \in L$ and $|y| \geq n$:

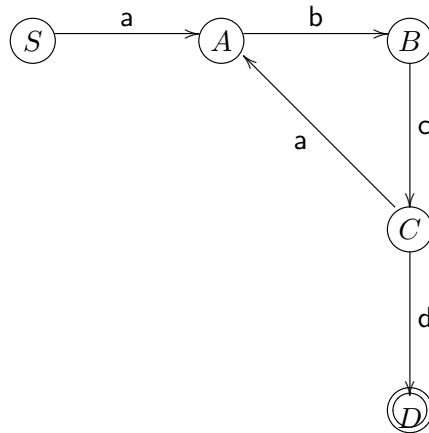
there exist u, v, w : $y = uvw$ and $|v| > 0$:

for all $i \in \mathbb{N}$: $xuv^i w \in L$

Note that $|y|$ denotes the length of the string y . Also remember that ‘for all $x \in X : \dots$ ’ is true if $X = \emptyset$, and ‘there exists $x \in X : \dots$ ’ is false if $X = \emptyset$.

9. Pumping Lemmas: the expressive power of languages

For example, consider the following automaton.



This automaton accepts: $abcabcd$, $abcabcabcd$, and, in general, $a(bca)^*bcd$. The statement of the pumping lemma amounts to the following. Take for n the number of states in the automaton (5). Let x, y, z be such that $xyz \in L$, and $|y| \geq n$. Then we know that in order to accept y , the above automaton has to pass at least twice through state A . The part that is accepted in between the two moments the automaton passes through state A can be pumped up to create sentences that contain an arbitrary number of copies of the string $v = bca$.

This pumping lemma is useful in showing that a language does *not* belong to the family of regular languages. Its application is typical of pumping lemmas in general; they are used *negatively* to show that a given language does not belong to some family.

Theorem 9.1 enables us to prove that a language L is not regular by showing that

for all $n \in \mathbb{N}$:
 there exist x, y, z : $xyz \in L$ and $|y| \geq n$:
 for all u, v, w : $y = uvw$ and $|v| > 0$:
 there exists $i \in \mathbb{N}$: $xuv^i wz \notin L$

In all applications of the pumping lemma in this chapter, this is the formulation we will use.

Note that if $n = 0$, we can choose $y = \epsilon$, and since there is no v with $|v| > 0$ such that $y = uvw$, the statement above holds for all such v (namely none!).

As an example, we will prove that language $L = \{a^m b^m \mid m \geq 0\}$ is not regular.

Let $n \in \mathbb{N}$.

Take $s = a^n b^n$ with $x = \epsilon$, $y = a^n$, and $z = b^n$.

Let u, v, w be such that $y = uvw$ with $v \neq \epsilon$, that is, $u = a^p$, $v = a^q$ and $w = a^r$ with $p + q + r = n$ and $q > 0$.

Take $i = 2$, then

$$\begin{aligned}
& xv^2wz \notin L \\
\Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
& a^{p+2q+r}b^n \notin L \\
\Leftarrow & p + q + r = n \\
& n + q \neq n \\
\Leftarrow & \text{ arithmetic} \\
& q > 0 \\
\Leftarrow & q > 0 \\
& \text{true}
\end{aligned}$$

Note that the language $L = \{a^m b^m \mid m \geq 0\}$ is context-free, and together with the fact that each regular grammar is also a context-free grammar it follows immediately that the set of regular languages is strictly smaller than the set of context-free languages.

Note that here we *use* the pumping lemma (and not the proof of the pumping lemma) to prove that a language is not regular. This kind of proof can be viewed as a kind of game: ‘for all’ is about an arbitrary element which can be chosen by the opponent; ‘there exists’ is about a particular element which you may choose. Choosing the right elements helps you ‘win’ the game, where winning means proving that a language is not regular.

Exercise 9.1. Prove that the following language is not regular

$$\{a^{k^2} \mid k \geq 0\}$$

Exercise 9.2. Show that the following language is not regular.

$$\{x \mid x \in \{a, b\}^* \wedge nr \text{ a } x < nr \text{ b } x\}$$

where $nr \text{ a } x$ is the number of occurrences of **a** in x .

Exercise 9.3. Prove that the following language is not regular

$$\{a^k b^m \mid k \leq m \leq 2k\}$$

Exercise 9.4. Show that the following language is not regular.

$$\{a^k b^l a^m \mid k > 5 \wedge l > 3 \wedge m \leq l\}$$

9.3. The pumping lemma for context-free languages

The Pumping Lemma for context-free languages gives a property that is satisfied by all context-free languages. This property is a statement of the form: in sentences exceeding a certain length, two sublists of bounded length can be identified that can be

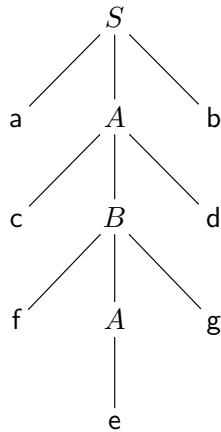
9. Pumping Lemmas: the expressive power of languages

duplicated while retaining a sentence. The idea behind this property is the following. Context-free languages are described by context-free grammars. For each sentence in the language there exists a derivation tree. When sentences have a derivation tree that is higher than the number of nonterminals, then at least one nonterminal will occur twice in a node; consequently a subtree can be inserted as often as desired.

As an example of an application of the Pumping Lemma, consider the context-free grammar with the following productions.

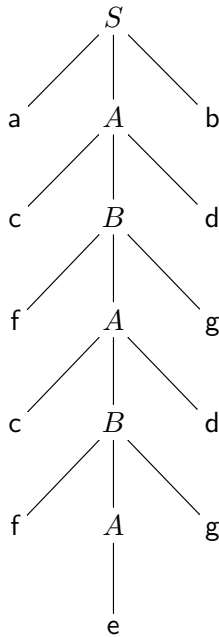
$$\begin{aligned} S &\rightarrow aAb \\ A &\rightarrow cBd \\ A &\rightarrow e \\ B &\rightarrow fAg \end{aligned}$$

The following parse tree represents the derivation of the sentence $acfegdb$.



If we replace the subtree rooted by the lower occurrence of nonterminal A by the

subtree rooted by the upper occurrence of A , we obtain the following parse tree.



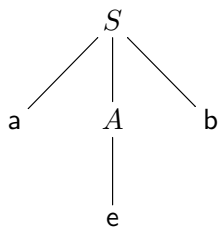
This parse tree represents the derivation of the sentence $acfcfegdgdb$. Thus we ‘pump’ the derivation of sentence $acfedgb$ to the derivation of sentence $acfcfegdgdb$. Repeating this step once more, we obtain a parse tree for the sentence

$acfcfcfegdgdb$

We can repeatedly apply this process to obtain derivation trees for all sentences of the form

$a(cf)^ie(gd)^ib$

for $i \geq 0$. The case $i = 0$ is obtained if we replace in the parse tree for the sentence $acfedgb$ the subtree rooted by the upper occurrence of nonterminal A by the subtree rooted by the lower occurrence of A :



This is a derivation tree for the sentence aeb . This step can be viewed as a negative pumping step.

The proof of the following lemma is given in Section 9.4.

9. Pumping Lemmas: the expressive power of languages

Theorem 9.2 (Context-free Pumping Lemma). *Let L be a context-free language. Then*

there exist c, d $: c, d \in \mathbb{N}$ $:$
for all z $: z \in L$ and $|z| \geq c$ $:$
there exist u, v, w, x, y $: z = uvwxy$ and $|vx| > 0$ and $|vwx| \leq d$ $:$
for all $i \in \mathbb{N}$ $: uv^iwx^iy \in L$

The Pumping Lemma is a tool with which we prove that a given language is not context-free. The proof obligation is to show that the property shared by all context-free languages does not hold for the language under consideration.

Theorem 9.2 enables us to prove that a language L is not context-free by showing that

for all c, d $: c, d \in \mathbb{N}$ $:$
there exists z $: z \in L$ and $|z| \geq c$ $:$
for all u, v, w, x, y $: z = uvwxy$ and $|vx| > 0$ and $|vwx| \leq d$ $:$
there exists $i \in \mathbb{N}$ $: uv^iwx^iy \notin L$

As an example, we will prove that the language T defined by

$$T = \{a^n b^n c^n \mid n > 0\}$$

is not context-free.

Proof. Let $c, d \in \mathbb{N}$.

Take $z = a^r b^r c^r$ with $r = \max(c, d)$.

Let u, v, w, x, y be such that $z = uvwxy$, $|vx| > 0$ and $|vwx| \leq d$

Note that our choice for r guarantees that substring vwx has one of the following shapes:

- vwx consists of just a's, or just b's, or just c's.
- vwx contains both a's and b's, or both b's and c's.

So vwx does *not* contain a's, b's, and c's.

Take $i = 0$, then

- If vwx consists of just a's, or just b's, or just c's, then it is impossible to write the string uvw as $a^s b^s c^s$ for some s , since only the number of terminals of one kind is decreased.

- If vw contains both a's and b's, or both b's and c's it lies somewhere on the border between a's and b's, or on the border between b's and c's. Then the string uvw can be written as

$$\begin{aligned}uvw &= a^s b^t c^r \\uvw &= a^r b^p c^q\end{aligned}$$

for some s, t, p, q , respectively. At least one of s and t or of p and q is less than r . Again this list is not an element of T .

□

Exercise 9.5. Why does vw not contain a's, b's, and c's?

Exercise 9.6. Prove that the following language is not context-free

$$\{a^{k^2} \mid k \geq 0\}$$

Exercise 9.7. Prove that the following language is not context-free

$$\{a^i \mid i \text{ is a prime number}\}$$

Exercise 9.8. Prove that the following language is not context-free

$$\{ww \mid w \in \{a, b\}^*\}$$

9.4. Proofs of pumping lemmas

This section gives the proof of the pumping lemmas.

9.4.1. Proof of the Regular Pumping Lemma, Theorem 9.1

Since L is a regular language, there exists a deterministic finite-state automaton D such that $L = L_{dfa} D$.

Take for n the number of states of D .

Let s be an element of L with sublist y such that $|y| \geq n$, say $s = xyz$.

Consider the sequence of states D passes through while processing y . Since $|y| \geq n$, this sequence has more than n entries, hence at least one state, say state A , occurs twice.

Take u, v, w as follows

- u is the initial part of y processed until the first occurrence of A ,
- v is the (nonempty) part of y processed from the first to the second occurrence of A ,
- w is the remaining part of y

9. Pumping Lemmas: the expressive power of languages

Note that D could have skipped processing v , and hence would have accepted $xuwz$. Furthermore, D can repeat the processing in between the first occurrence of A and the second occurrence of A as often as desired, and hence it accepts $xuv^i w z$ for all $i \in \mathbb{N}$. Formally, a simple proof by induction shows that $(\forall i : i \geq 0 : xuv^i w z \in L)$.

9.4.2. Proof of the Context-free Pumping Lemma, Theorem 9.2

Let $G = (T, N, R, S)$ be a context-free grammar such that $L = L(G)$. Let m be the length of the longest right-hand side of any production, and let k be the number of nonterminals of G .

Take $c = m^k$. In Lemma 9.3 below we prove that if z is a list with $|z| > c$, then in all derivation trees for z there exists a path of length at least $k+1$.

Let $z \in L$ such that $|z| > c$. Since grammar G has k nonterminals, there is at least one nonterminal that occurs more than once in a path of length $k+1$ (which contains $k+2$ symbols, of which at most one is a terminal, and all others are nonterminals) of a derivation tree for z . Consider the nonterminal A that satisfies the following requirements.

- A occurs at least twice in the path of length $k+1$ of a derivation tree for z .

Call the list corresponding to the derivation tree rooted at the lower A w , and call the list corresponding to the derivation tree rooted at the upper A (which contains the list w) vwx .

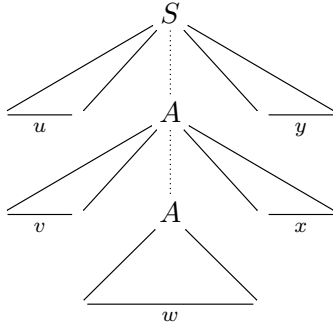
- A is chosen such that at most one of v and x equals ϵ .
- Finally, we suppose that below the upper occurrence of A no other nonterminal that satisfies the same requirements occurs, that is, the two A 's are the lowest pair of nonterminals satisfying the above requirements.

First we show that a nonterminal A satisfying the above requirements exists. We prove this statement by contradiction. Suppose that for all nonterminals A that occur twice on a path of length at least $k+1$ both v and x , the border lists of the list vwx corresponding to the tree rooted at the upper occurrence of A , are both equal to ϵ . Then we can replace the tree rooted at the upper A by the tree rooted at the lower A without changing the list corresponding to the tree. Thus we can replace all paths of length at least $k+1$ by a path of length at most k . But this contradicts Lemma 9.3 below, and we have a contradiction. It follows that a nonterminal A satisfying the above requirements exists.

There exists a derivation tree for z in which the path from the upper A to the leaf has length at most $k+1$, since either below A no nonterminal occurs twice, or there is one or more nonterminal B that occurs twice, but the border lists v' and x' from the list $v'w'x'$ corresponding to the tree rooted at the upper occurrence of nonterminal B are empty. Since the lists v' and x' are empty, we can replace the tree rooted at

the upper occurrence of B by the tree rooted at the lower occurrence of B without changing the list corresponding to the derivation tree. Since we can do this for all nonterminals B that occur twice below the upper occurrence of A , there exists a derivation tree for z in which the path from the upper A to the leaf has length at most $k+1$. It follows from Lemma 9.3 below that the length of vwx is at most m^{k+1} , so we define $d = m^{k+1}$.

Suppose $z = uvwxy$, that is, the list corresponding to the subtree to the left (right) of the upper occurrence of A is u (y). This situation is depicted as follows.



We prove by induction that $(\forall i : i \geq 0 : uv^iwx^i y \in L)$. In this proof we apply the tree substitution process described in the example before the lemma. For $i = 0$ we have to show that the list $uvwxy$ is a sentence in L . The list $uvwxy$ is obtained if the tree rooted at the upper A is replaced by the tree rooted at the lower A . Suppose that for all $i \leq n$ we have $uv^iwx^i y \in L$. The list $uv^{i+1}wx^{i+1}y$ is obtained if the tree rooted at the lower A in the derivation tree for $uv^iwx^i y \in L$ is replaced by the tree rooted above it A . This proves the induction step.

The proof of the Pumping Lemma above frequently refers to the following lemma.

Theorem 9.3. *Let G be a context-free grammar, and suppose that the longest right-hand side of any production has length m . Let t be a derivation tree for a sentence $z \in L(G)$. If $\text{height } t \leq j$, then $|z| \leq m^j$.*

Proof. We prove a slightly stronger result: if t is a derivation tree for a list z , but the root of t is not necessarily the start-symbol, and $\text{height } t \leq j$, then $|z| \leq m^j$. We prove this statement by induction on j .

For the base case, suppose $j = 1$. Then tree t corresponds to a single production in G , and since the longest right-hand side of any production has length m , we have that $|z| \leq m = m^j$.

For the induction step, assume that for all derivation trees t of height at most j we have that $|z| \leq m^j$, where z is the list corresponding to t . Suppose we have a tree t of height $j+1$. Let A be the root of t , and suppose the top of the tree

9. Pumping Lemmas: the expressive power of languages

corresponds to the production $A \rightarrow v$ in G . For all trees s rooted at the symbols of v we have $\text{height } s \leq j$, so the induction hypothesis applies to these trees, and the lists corresponding to the trees rooted at the symbols of v all have length at most m^j . Since $A \rightarrow v$ is a production of G , and since the longest right-hand side of any production has length m , the list corresponding to the tree rooted at A has length at most $m \times m^j = m^{j+1}$, which proves the induction step. \square

Summary

This section introduces pumping lemmas. Pumping lemmas are used to prove that languages are not regular or not context-free.

9.5. Exercises

Exercise 9.9. Show that the following language is not regular.

$$\{x \mid x \in \{0,1\}^* \wedge nr \ 1 \ x = nr \ 0 \ x\}$$

where function nr takes an element a and a list x , and returns the number of occurrences of a in x .

Exercise 9.10. Consider the following language:

$$\{a^i b^j \mid 0 \leq i \leq j\}$$

1. Is this language context-free? If it is, give a context-free grammar and prove that this grammar generates the language. If it is not, why not?
2. Is this language regular? If it is, give a regular grammar and prove that this grammar generates the language. If it is not, why not?

Exercise 9.11. Consider the following language:

$$\{w c w \mid w \in \{a, b\}^*\}$$

1. Is this language context-free? If it is, give a context-free grammar and prove that this grammar generates the language. If it is not, why not?
2. Is this language regular? If it is, give a regular grammar and prove that this grammar generates the language. If it is not, why not?

Exercise 9.12. Consider the grammar G with the following productions.

$$\left\{ \begin{array}{l} S \rightarrow \{A\} \\ S \rightarrow \epsilon \\ A \rightarrow S \\ A \rightarrow AA \\ A \rightarrow \}\{ \end{array} \right.$$

1. Is this grammar

- Context-free?
- Regular?

Why?

2. Give the language of G without referring to G itself. Prove that your description is correct.
3. Is the language of G
 - Context-free?
 - Regular?

Why?

Exercise 9.13. Consider the grammar G with the following productions.

$$\left\{ \begin{array}{l} S \rightarrow \epsilon \\ S \rightarrow 0 \\ S \rightarrow 1 \\ S \rightarrow S0 \end{array} \right.$$

1. Is this grammar
 - Context-free?
 - Regular?

Why?

2. Give the language of G without referring to G itself. Prove that your description is correct.
3. Is the language of G
 - Context-free?
 - Regular?

Why?

9. Pumping Lemmas: the expressive power of languages

10. LL Parsing

Introduction

This chapter introduces LL(1) parsing. LL(1) parsing is an efficient (linear in the length of the input string) method for parsing that can be used for all LL(1) grammars. A grammar is LL(1) if at each step in a derivation the next symbol in the input uniquely determines the production that should be applied. In order to determine whether or not a grammar is LL(1), we introduce several kinds of grammar analyses, such as determining whether or not a nonterminal can derive the empty string, and determining the set of symbols that can appear as the first symbol in a derivation from a nonterminal.

Goals

After studying this chapter you will

- know the definition of LL(1) grammars;
- know how to parse a sentence from an LL(1) grammar;
- be able to apply different kinds of grammar analyses in order to determine whether or not a grammar is LL(1).

This chapter is organised as follows. Section 10.1 describes the background of LL(1) parsing, and Section 10.2 describes an implementation in Haskell of LL(1) parsing and the different kinds of grammar analyses needed for checking whether or not a grammar is LL(1).

10.1. LL Parsing: Background

In the previous chapters we have shown how to construct parsers for sentences of context-free languages using combinator parsers. Since these parsers may backtrack, the resulting parsers are sometimes a bit slow. There are several ways in which we can put extra restrictions on context-free grammars such that we can parse sentences of the corresponding languages efficiently. This chapter discusses one such restriction: *LL*(1). Other restrictions, not discussed in these lecture notes are *LR*(1), *LALR*(1), *SLR*(1), etc.

10.1.1. A stack machine for parsing

This section presents a stack machine for parsing sentences of context-free grammars. We will use this machine in the following subsections to illustrate why we need grammar analysis.

The stack machine we use in this section differs from the stack machines introduced in Sections 5.4.4 and 7.2. A stack machine for a grammar G has a stack and an input, and performs one of the following two actions.

1. **Expand:** If the top stack symbol is a nonterminal, it is popped from the stack and a right-hand side from a production of G for the nonterminal is pushed onto the stack. The production is chosen nondeterministically.
2. **Match:** If the top stack symbol is a terminal, then it is popped from the stack and compared with the next symbol of the input sequence. If they are equal, then this terminal symbol is ‘read’. If the stack symbol and the next input symbol do not match, the machine signals an error, and the input sentence cannot be accepted.

These actions are performed until the stack is empty. A stack machine for G accepts an input if it can terminate with an empty input when starting with the start-symbol from G on the stack.

For example, let grammar G be the grammar with the productions:

$$S \rightarrow aS \mid cS \mid b$$

The stack machine for G accepts the input string aab because it can perform the following actions (the first component of the state (before the | in the picture below) is the symbol stack and the second component of the state (after the |) is the unmatched (remaining part of the) input string):

stack	input
S	aab
aS	aab
S	ab
aS	ab
S	b
b	b

and end with an empty input. However, if the machine had chosen the production $S \rightarrow cS$ in its first step, it would have been stuck. So not all possible sequences of actions from the *state* (S, aab) lead to an empty input. If there is at least one sequence of actions that ends with an empty input on an input string, the input

string is accepted. In this sense, the stack machine is similar to a nondeterministic finite-state automaton.

10.1.2. Some example derivations

This section gives three examples in which the stack machine for parsing is applied. It turns out that, for all three examples, the nondeterministic stack machine can act in a deterministic way by looking ahead one (or two) symbols of the sequence of input symbols. Each of the examples exemplifies why different kinds of grammar analyses are useful in parsing.

The first example

Our first example is **gramm1**. The set of terminal symbols of **gramm1** is $\{a, b, c\}$, the set of nonterminal symbols is $\{S, A, B, C\}$, the start symbol is S ; and the productions are

$$\begin{aligned} S &\rightarrow cA \mid b \\ A &\rightarrow cBC \mid bSA \mid a \\ B &\rightarrow cc \mid Cb \\ C &\rightarrow aS \mid ba \end{aligned}$$

We want to know whether or not the string **ccccba** is a sentence of the language of this grammar. The stack machine produces, amongst others, the following sequence, corresponding with a leftmost derivation of **ccccba**.

stack	input
S	ccccba
cA	ccccba
A	cccba
cBC	cccba
BC	ccba
ccC	ccba
cC	cba
C	ba
ba	ba
a	a

Starting with S the machine chooses between two productions:

$$S \rightarrow cA \mid b$$

but, since the first symbol of the string `ccccba` to be recognised is `c`, the only applicable production is the first one. After expanding `S` a match-action removes the leading `c` from `ccccba` and `cA`. So now we have to derive the string `cccba` from `A`. The machine chooses between three productions:

$$A \rightarrow cBC \mid bSA \mid a$$

and again, since the next symbol of the remaining string `cccba` to be recognised is `c`, the only applicable production is the first one. After expanding `A` a match-action removes the leading `c` from `cccba` and `cBC`. So now we have to derive the string `ccba` from `BC`. The top stack symbol of `BC` is `B`. The machine chooses between two productions:

$$B \rightarrow cc \mid Cb$$

The next symbol of the remaining string `ccba` to be recognised is, once again, `c`. The first production is applicable, but the second production may be applicable as well. To decide whether it also applies we have to determine the symbols that can appear as the first element of a string derived from `B` starting with the second production. The first symbol of the alternative `Cb` is the nonterminal `C`. From the productions

$$C \rightarrow aS \mid ba$$

it is immediately clear that a string derived from `C` starts with either an `a` or a `b`. The set $\{a, b\}$ is called the *first* set of the nonterminal `C`. Since the next symbol in the remaining string to be recognised is a `c`, the second production cannot be applied. After expanding `B` and performing two match-actions it remains to derive the string `ba` from `C`. The machine chooses between two productions $C \rightarrow aS$ and $C \rightarrow ba$. Clearly, only the second one applies, and, after two match-actions, leads to success.

From the above derivation we conclude the following.

Deriving the sentence `ccccba` using `gramm1` is a deterministic computation: at each step of the derivation there is only one applicable alternative for the nonterminal on top of the stack.

Determinicity is obtained by looking at the set of firsts of the nonterminals.

The second example

A second example is the grammar `gramm2` whose productions are

$$\begin{aligned} S &\rightarrow abA \mid aa \\ A &\rightarrow bb \mid bS \end{aligned}$$

Now we want to know whether or not the string `abbb` is a sentence of the language of this grammar. The stack machine produces, amongst others, the following sequence

stack	input
S	abbb
abA	abbb
bA	bbb
A	bb
bb	bb
b	b

Starting with S the machine chooses between two productions:

$$S \rightarrow abA \mid aa$$

since both alternatives start with an a , it is not sufficient to look at the first symbol a of the string to be recognised. The problem is that the *lookahead sets* (the *lookahead set* of a production $N \rightarrow \alpha$ is the set of terminal symbols that can appear as the first symbol of a string that can be derived from N starting with the production $N \rightarrow \alpha$, the definition is given in the following subsection) of the two productions for S both contain a . However, if we look at the first two symbols ab , then we find that the only applicable production is the first one. After expanding and matching it remains to derive the string bb from A . Again, looking ahead one symbol in the input string does not give sufficient information for choosing one of the two productions

$$A \rightarrow bb \mid bS$$

for A . If we look at the first two symbols bb of the input string, then we find that the first production applies (and, after matching, leads to success). Each string derived from A starting with the second production starts with a b and, since it is not possible to derive a string starting with another b from S , the second production does not apply.

From the above derivation we conclude the following.

Deriving the string **abbb** using **gramm2** is a deterministic computation: at each step of the derivation there is only one applicable alternative for the nonterminal on the top of the stack.

Again, determinicity is obtained by analysing the set of firsts (of strings of length 2) of the nonterminals. Alternatively, we can left-factor the grammar to obtain a grammar in which all productions for a nonterminal start with a different terminal symbol.

The third example

A third example is grammar `gramm3` with the following productions:

$$\begin{aligned} S &\rightarrow AaS \mid B \\ A &\rightarrow cS \mid \epsilon \\ B &\rightarrow \mathbf{b} \end{aligned}$$

Now we want to know whether or not the string `acbabb` is an element of the language of this grammar. The stack machine produces the following sequence

stack	input
S	acbabb
AaS	acbabb
aS	acbabb

Starting with S the machine chooses between two productions:

$$S \rightarrow AaS \mid B$$

since each nonempty string derived from A starts with a `c`, and each nonempty string derived from B starts with a `b`, there does not seem to be a candidate production to start a leftmost derivation of `acbabb` with. However, since A can also derive the empty string, we can apply the first production, and then apply the empty string for A , producing `aS` which, as required, starts with an `a`. We do not explain the rest of the leftmost derivation since it does not use any empty strings any more.

Nonterminal symbols that can derive the empty sequence will play a central role in the grammar analysis problems which we will consider in Section 10.2.

From the above derivation we conclude the following.

Deriving the string `acbabb` using `gramm3` is a deterministic computation: at each step of the derivation there is only one applicable alternative for the nonterminal on the top of the stack.

Determinicity is obtained by analysing whether or not nonterminals can derive the empty string, and which terminal symbols can follow upon a nonterminal in a derivation.

10.1.3. $LL(1)$ grammars

The examples in the previous subsection show that the derivations of the example sentences are deterministic, provided we can look ahead one or two symbols in the input. An obvious question now is: for which grammars are all derivations deterministic? Of course, as the second example shows, the answer to this question depends on the number of symbols we are allowed to look ahead. In the rest of this chapter we assume that we may look 1 symbol ahead. A grammar for which all derivations are deterministic with 1 symbol lookahead is called $LL(1)$: Leftmost with a Lookahead of 1. Since all derivations of sentences of $LL(1)$ grammars are deterministic, $LL(1)$ is a desirable property of grammars.

To formalise this definition, we define `lookAhead` sets.

Definition 10.1 (`lookAhead` set). The `lookahead set` of a production $N \rightarrow \alpha$ is the set of terminal symbols that can appear as the first symbol of a string that can be derived from $N\delta$ (where $N\delta$ appears as a tail substring in a derivation from the start-symbol) starting with the production $N \rightarrow \alpha$. So

$$\text{lookAhead}(N \rightarrow \alpha) = \{x \mid S \xRightarrow{*} \gamma N \delta \Rightarrow \gamma \alpha \delta \xRightarrow{*} \gamma x \beta\}$$

For example, for the productions of `gramm1` we have

$$\begin{aligned} \text{lookAhead}(S \rightarrow cA) &= \{c\} \\ \text{lookAhead}(S \rightarrow b) &= \{b\} \\ \text{lookAhead}(A \rightarrow cBC) &= \{c\} \\ \text{lookAhead}(A \rightarrow bSA) &= \{b\} \\ \text{lookAhead}(A \rightarrow a) &= \{a\} \\ \text{lookAhead}(B \rightarrow cc) &= \{c\} \\ \text{lookAhead}(B \rightarrow C'b) &= \{a, b\} \\ \text{lookAhead}(C \rightarrow aS) &= \{a\} \\ \text{lookAhead}(C \rightarrow ba) &= \{b\} \end{aligned}$$

We use `lookAhead` sets in the definition of $LL(1)$ grammar.

Definition 10.2 ($LL(1)$ grammar). A grammar G is $LL(1)$ if all pairs of productions of the same nonterminal have disjoint lookahead sets, that is: for all productions $N \rightarrow \alpha, N \rightarrow \beta$ of G :

$$\text{lookAhead}(N \rightarrow \alpha) \cap \text{lookAhead}(N \rightarrow \beta) = \emptyset$$

Since all `lookAhead` sets for productions of the same nonterminal of `gramm1` are disjoint, `gramm1` is an $LL(1)$ grammar. For `gramm2` we have:

$$\begin{aligned} \text{lookAhead}(S \rightarrow abA) &= \{a\} \\ \text{lookAhead}(S \rightarrow aa) &= \{a\} \\ \text{lookAhead}(A \rightarrow bb) &= \{b\} \\ \text{lookAhead}(A \rightarrow bS) &= \{b\} \end{aligned}$$

Here, the `lookAhead` sets for both nonterminals S and A are not disjoint, and it follows that `gramm2` is not LL(1). `gramm2` is an LL(2) grammar, where an LL(k) grammar for $k \geq 2$ is defined similarly to an LL(1) grammar: instead of one symbol lookahead we have k symbols lookahead.

How do we determine whether or not a grammar is LL(1)? Clearly, to answer this question we need to know the lookahead sets of the productions of the grammar. The `lookAhead` set of a production $N \rightarrow \alpha$, where α starts with a terminal symbol x , is simply x . But what if α starts with a nonterminal P , that is $\alpha = P\beta$, for some β ? Then we have to determine the set of terminal symbols with which strings derived from P can start. But if P can derive the empty string, we also have to determine the set of terminal symbols with which a string derived from β can start. As you see, in order to determine the `lookAhead` sets of productions, we are interested in

- whether or not a nonterminal can derive the empty string (`empty`);
- which terminal symbols can appear as the first symbol in a string derived from a nonterminal (`firsts`);
- and which terminal symbols can follow upon a nonterminal in a derivation (`follow`).

In each of the following definitions we assume that a grammar G is given.

Definition 10.3 (Empty). Function `empty` takes a nonterminal N , and determines whether or not the empty string can be derived from the nonterminal:

$$\text{empty } N = N \xRightarrow{*} \epsilon$$

For example, for `gramm3` we have:

$$\begin{aligned} \text{empty } S &= \text{False} \\ \text{empty } A &= \text{True} \\ \text{empty } B &= \text{False} \end{aligned}$$

Definition 10.4 (First). The *set of firsts* of a nonterminal N is the set of terminal symbols that can appear as the first symbol of a string that can be derived from N :

$$\text{firsts } N = \{x \mid N \xRightarrow{*} x\beta\}$$

For example, for `gramm3` we have:

$$\begin{aligned} \text{firsts } S &= \{\mathbf{a}, \mathbf{b}, \mathbf{c}\} \\ \text{firsts } A &= \{\mathbf{c}\} \\ \text{firsts } B &= \{\mathbf{b}\} \end{aligned}$$

We could have given more restricted definitions of `empty` and `firsts`, by only looking at derivations from the start-symbol, for example,

$$\text{empty } N = S \xRightarrow{*} \alpha N \beta \xRightarrow{*} \alpha \beta$$

but the simpler definition above suffices for our purposes.

Definition 10.5 (Follow). The *follow set* of a nonterminal N is the set of terminal symbols that can follow on N in a derivation starting with the start-symbol S from the grammar G :

$$\text{follow } N = \{x \mid S \xRightarrow{*} \alpha N x \beta\}$$

For example, for `gramm3` we have:

$$\begin{aligned} \text{follow } S &= \{a\} \\ \text{follow } A &= \{a\} \\ \text{follow } B &= \{a\} \end{aligned}$$

In the following section we will give programs with which `lookahead`, `empty`, `firsts`, and `follow` are computed.

Exercise 10.1. Give the results of the function `empty` for the grammars `gramm1` and `gramm2`.

Exercise 10.2. Give the results of the function `firsts` for the grammars `gramm1` and `gramm2`.

Exercise 10.3. Give the results of the function `follow` for the grammars `gramm1` and `gramm2`.

Exercise 10.4. Give the results of the function `lookahead` for grammar `gramm3`. Is `gramm3` an LL(1) grammar ?

Exercise 10.5. Grammar `gramm2` is not LL(1), but it can be transformed into an LL(1) grammar by left factoring. Give this equivalent grammar `gramm2'` and give the results of the functions `empty`, `first`, `follow` and `lookAhead` on this grammar. Is `gramm2'` an LL(1) grammar?

Exercise 10.6. A non-leftrecursive grammar for `Bit-Lists` is given by the following grammar (see your answer to Exercise 2.18):

$$\begin{aligned} L &\rightarrow BR \\ R &\rightarrow \epsilon \mid ,BR \\ B &\rightarrow 0 \mid 1 \end{aligned}$$

Give the results of functions `empty`, `firsts`, `follow` and `lookAhead` on this grammar. Is this grammar LL(1)?

10.2. LL Parsing: Implementation

Until now we have written parsers with parser combinators. Parser combinators use backtracking, and this is sometimes a cause of inefficiency. If a grammar is LL(1) we do not need backtracking anymore: parsing is deterministic. We can use this fact by either adjusting the parser combinators so that they don't use backtracking anymore, or by writing a special purpose LL(1) parsing program. We present the latter in this section.

This section describes the implementation of a program that parses sentences of LL(1) grammars. The program works for arbitrary context-free LL(1) grammars, so we first describe how to represent context-free grammars in Haskell. Another consequence of the fact that the program parses sentences of arbitrary context-free LL(1) grammars is that we need a generic representation of parse trees in Haskell. The second subsection defines a datatype for parse trees in Haskell. The third subsection presents the program that parses sentences of LL(1) grammars. This program assumes that the input grammar is LL(1), so in the fourth subsection we give a function that determines whether or not a grammar is LL(1). Both this and the LL(1) parsing function use a function that determines the lookahead of a production. This function is presented in the fifth subsection. The last subsections of this section define functions for determining the empty, first, and follow symbols of a nonterminal.

10.2.1. Context-free grammars in Haskell

A context-free grammar may be represented by a pair: its start symbol, and its productions. How do we represent terminal and nonterminal symbols? There are at least two possibilities.

- The rigorous approach uses a datatype `Symbol`:

```
data Symbol a b = N a | T b
```

The advantage of this approach is that nonterminals and terminals are strictly separated, the disadvantage is the notational overhead of constructors that has to be carried around. However, a rigorous implementation of context-free grammars should keep terminals and nonterminals apart, so this is the preferred implementation. But in this section we will use the following implementation:

- ```
class Eq s => Symbol s where
 isT :: s -> Bool
 isN :: s -> Bool
 isT = not . isN
```

where `isN` and `isT` determine whether or not a symbol is a nonterminal or a terminal, respectively. This notation is compact, but terminals and nonterminals are no longer strictly separated, and symbols that are used as nonterminals

cannot be used as terminals anymore. For example, the type of characters is made an instance of this class by defining:

```
instance Symbol Char where
 isN c = 'A' <= c && c <= 'Z'
```

that is, capitals are nonterminals, and, by definition, all other characters are terminals.

A context-free grammar is a value of the type CFG:

```
type CFG s = (s, [(s, [s])])
```

where the list in the second component associates nonterminals to right-hand sides. So an element of this list is a production. For example, the grammar with productions

$$\begin{aligned} S &\rightarrow AaS \mid B \mid CB \\ A &\rightarrow SC \mid \epsilon \\ B &\rightarrow A \mid b \\ C &\rightarrow D \\ D &\rightarrow d \end{aligned}$$

is represented as:

```
exGrammar :: CFG Char
exGrammar =
 ('S', [(('S', "AaS"), ('S', "B"), ('S', "CB"))
 , ('A', "SC"), ('A', "")
 , ('B', "A"), ('B', "b")
 , ('C', "D")
 , ('D', "d")
]
)
```

On this type we define some functions for extracting the productions, nonterminals, terminals, etc. from a grammar.

```
start :: CFG s -> s
start = fst

prods :: CFG s -> [(s, [s])]
prods = snd

terminals :: (Symbol s, Ord s) => CFG s -> [s]
terminals = unions . map (filter isT . snd) . snd
```

## 10. LL Parsing

where `unions :: Ord s => [[s]] -> [s]` returns the union of the ‘sets’ of symbols in the lists.

```
nonterminals :: (Symbol s, Ord s) => CFG s -> [s]
nonterminals = nub . map fst . snd
```

Here, `nub :: Ord s => [s] -> [s]` removes duplicates from a list.

```
symbols :: (Symbol s, Ord s) => CFG s -> [s]
symbols grammar =
 union (terminals grammar) (nonterminals grammar)

nt2prods :: Eq s => CFG s -> s -> [(s,[s])]
nt2prods grammar s =
 filter (\(nt,rhs) -> nt==s) (prods grammar)
```

where function `union` returns the set union of two lists (removing duplicates). For example, we have

```
?start exGrammar
'S'

? terminals exGrammar
"abd"
```

### 10.2.2. Parse trees in Haskell

A parse tree is a tree, in which each internal node is labelled with a nonterminal, and has a list of children (corresponding with a right-hand side of a production of the nonterminal). It follows that parse trees can be represented as rose trees with symbols, where the datatype of rose trees is defined by:

```
data Rose a = Node a [Rose a] | Nil
```

The constructor `Nil` has been added to simplify error handling when parsing: when a sentence cannot be parsed, the ‘parse tree’ `Nil` is returned. Strictly speaking it should not occur in the datatype of rose trees.

### 10.2.3. LL(1) parsing

This section defines a function `ll1` that takes a grammar and a terminal string as input, and returns one tuple: a parse tree, and the rest of the inputstring that has not been parsed. So `ll1` takes a grammar, and returns a parser with `Rose s` as its result type. As mentioned in the beginning of this section, our parser doesn't need backtracking anymore, since parsing with an LL(1) grammar is deterministic. Therefore, the parser type is adjusted as follows:

```
type Parser b a = [b] -> (a, [b])
```

Using this parser type, the type of the function `ll1` is:

```
ll1 :: (Symbol s, Ord s) => CFG s -> Parser s (Rose s)
```

Function `ll1` is defined in terms of two functions. Function `isll1 :: CFG s -> Bool` is a function that checks whether or not a grammar is LL(1). And function `gll1` (for *generalised* LL(1)) produces a *list* of rose trees for a *list* of symbols. `ll1` is obtained from `gll1` by giving `gll1` the singleton list containing the start symbol of the grammar as argument.

```
ll1 grammar input =
 if isll1 grammar
 then let ([rose], rest) = gll1 grammar [start grammar] input
 in (rose, rest)
 else error "ll1: grammar not LL(1)"
```

So now we have to implement functions `isll1` and `gll1`. Function `isll1` is implemented in the following subsection. Function `gll1` also uses two functions. Function `grammar2ll1table` takes a grammar and returns the LL(1) table: the association list that associates productions with their lookahead sets. And function `choose` takes a terminal symbol, and chooses a production based on the LL(1) table.

```
gll1 :: (Symbol s, Ord s) => CFG s -> [s] -> Parser s [Rose s]
gll1 grammar =
 let ll1table = grammar2ll1table grammar
 -- The LL(1) table.
 nt2prods nt = filter (\((n,l),r) -> n==nt) ll1table
 -- nt2prods returns the productions for nonterminal
 -- nt from the LL(1) table
 selectprod nt t = choose t (nt2prods nt)
 -- selectprod selects the production for nt from the
```

## 10. LL Parsing

```

-- LL(1) table that should be taken when t is the next
-- symbol in the input.
in \stack input ->
 case stack of
 [] -> ([], input)
 (s:ss) ->
 if isT s
 then -- match
 let (rts,rest) = gll1 grammar ss (tail input)
 in if s == head input
 then (Node s []: rts, rest)
 -- The parse tree is a leaf (a node with
 -- no children).
 else ([Nil], input)
 -- The input cannot be parsed
 else -- expand
 let t = head input
 (rts,zs) = gll1 grammar (selectprod s t) input
 -- Try to parse according to the production
 -- obtained from the LL(1) table from s.
 (rrs,vs) = gll1 grammar ss zs
 -- Parse the rest of the symbols on the
 -- stack.
 in ((Node s rts): rrs, vs)

```

Functions `grammar2ll1table` and `choose`, which are used in the above function `gll1`, are defined as follows. These functions use function `lookaheadp`, which returns the lookahead set of a production and is defined in one of the following subsections.

```

grammar2ll1table :: (Symbol s, Ord s) => CFG s -> [((s,[s]),[s])]
grammar2ll1table grammar =
 map (\x -> (x,lookaheadp grammar x)) (prods grammar)

choose :: Eq a => a -> [((b,c),[a])] -> c
choose t l =
 let [((s,rhs), ys)] = filter (\((x,p),q) -> t `elem` q) l
 in rhs

```

Note that function `choose` assumes that there is exactly one element in the association list in which the element `t` occurs.

### 10.2.4. Implementation of isLL(1)

Function `isll1` checks whether or not a context-free grammar is LL(1). It does this by computing for each nonterminal of its argument grammar the set of lookahead sets (one set for each production of the nonterminal), and checking that all of these sets are disjoint. It is defined in terms of a function `lookaheadn`, which computes the lookahead sets of a nonterminal, and a function `disjoint`, which determines whether or not all sets in a list of sets are disjoint. All sets in a list of sets are disjoint if the length of the concatenation of these sets equals the length of the union of these sets.

```
isll1 :: (Symbol s, Ord s) => CFG s -> Bool
isll1 grammar =
 and (map (disjoint . lookaheadn grammar) (nonterminals grammar))

disjoint :: Ord s => [[s]] -> Bool
disjoint xss = length (concat xss) == length (unions xss)
```

Function `lookaheadn` computes the lookahead sets of a nonterminal by computing all productions of a nonterminal, and computing the lookahead set of each of these productions by means of function `lookaheadp`.

```
lookaheadn :: (Symbol s, Ord s) => CFG s -> s -> [[s]]
lookaheadn grammar =
 map (lookaheadp grammar) . nt2prods grammar
```

### 10.2.5. Implementation of lookahead

Function `lookaheadp` takes a grammar and a production, and returns the lookahead set of the production. It is defined in terms of four functions. Each of the first three functions will be defined in a separate subsection below, the fourth function is defined in this subsection.

- `isEmpty :: (Ord s, Symbol s) => CFG s -> s -> Bool`

Function `isEmpty` takes a grammar and a nonterminal and determines whether or not the empty string can be derived from the nonterminal in the grammar. (This function was called `empty` in Definition 10.3.)

- `firsts :: (Ord s, Symbol s) => CFG s -> [(s, [s])]`

Function `firsts` takes a grammar and computes the set of firsts of each symbol (the set of firsts of a terminal is the terminal itself).

## 10. LL Parsing

- `follow :: (Ord s, Symbol s) => CFG s -> [(s,[s])]`

Function `follow` takes a grammar and computes the follow set of each nonterminal (so it associates a list of symbols with each nonterminal).

- `lookSet :: Ord s =>`  
`(s -> Bool) -> -- isEmpty`  
`(s -> [s]) -> -- firsts?`  
`(s -> [s]) -> -- follow?`  
`(s, [s]) -> -- production`  
`[s] -- lookahead set`

Note that we use the operator `?`, see Section 5.4.2, on the `firsts` and `follow` association lists. Function `lookSet` takes a predicate, two functions that given a nonterminal return the first and follow set, respectively, and a production, and returns the lookahead set of the production. Function `lookSet` is introduced after the definition of function `lookaheadp`.

Now we define:

```
lookaheadp :: (Symbol s, Ord s) => CFG s -> (s,[s]) -> [s]
lookaheadp grammar =
 lookSet (isEmpty grammar) ((firsts grammar)?) ((follow grammar)?)
```

We will exemplify the definition of function `lookSet` with the grammar `exGrammar`, with the following productions:

$$\begin{aligned} S &\rightarrow AaS \mid B \mid CB \\ A &\rightarrow SC \mid \epsilon \\ B &\rightarrow A \mid \mathbf{b} \\ C &\rightarrow D \\ D &\rightarrow \mathbf{d} \end{aligned}$$

Consider the production  $S \rightarrow AaS$ . The lookahead set of the production contains the set of symbols which can appear as the first terminal symbol of a sequence of symbols derived from  $A$ . But, since the nonterminal symbol  $A$  can derive the empty string, the lookahead set also contains the symbol  $\mathbf{a}$ .

Consider the production  $A \rightarrow SC$ . The lookahead set of the production contains the set of symbols which can appear as the first terminal symbol of a sequence of symbols derived from  $S$ . But, since the nonterminal symbol  $S$  can derive the empty string, the lookahead set also contains the set of symbols which can appear as the first terminal symbol of a sequence of symbols derived from  $C$ .

Finally, consider the production  $B \rightarrow A$ . The lookahead set of the production contains the set of symbols which can appear as the first terminal symbol of a sequence



of symbols derived from  $A$ . But, since the nonterminal symbol  $A$  can derive the empty string, the lookahead set also contains the set of terminal symbols which can follow the nonterminal symbol  $B$  in some derivation.

The examples show that it is useful to have functions `firsts` and `follow` in which, for every nonterminal symbol  $n$ , we can look up the terminal symbols which can appear as the first terminal symbol of a sequence of symbols in some derivation from  $n$  and the set of terminal symbols which can follow the nonterminal symbol  $n$  in a sequence of symbols occurring in some derivation respectively. It turns out that the definition of function `follow` also makes use of a function `lasts` which is similar to the function `firsts`, but which deals with last nonterminal symbols rather than first terminal ones.

The examples also illustrate a control structure which will be used very often in the following algorithms: we will fold over right-hand sides. While doing so we compute sets of symbols for all the symbols of the right-hand side which we encounter and collect them into a final set of symbols. Whenever such a list for a symbol is computed, there are always two possibilities:

- either we continue folding and return the result of taking the union of the set obtained from the current element and the set obtained by recursively folding over the rest of the right-hand side
- or we stop folding and immediately return the set obtained from the current element.

We continue if the current symbol is a nonterminal which can derive the empty sequence and we stop if the current symbol is either a terminal symbol or a nonterminal symbol which cannot derive the empty sequence. The following function makes this statement more precise.

```
foldrRhs :: Ord s =>
 (s -> Bool) ->
 (s -> [s]) ->
 [s] ->
 [s] ->
 [s]
foldrRhs p f start = foldr op start
 where op x xs = f x 'union' if p x then xs else []
```

The function `foldrRhs` is, of course, most naturally defined in terms of the function `foldr`. This function is somewhere in between a general purpose and an application specific function (we could easily have made it more general though). In the exercises we give an alternative characterisation of `foldRhs`. We will also need a function `scanrRhs` which is like `foldrRhs` but accumulates intermediate results in a list. The function `scanrRhs` is most naturally defined in terms of the function `scanr`.

## 10. LL Parsing

```

scanrRhs :: Ord s =>
 (s -> Bool) ->
 (s -> [s]) ->
 [s] ->
 [s] ->
 [[s]]
scanrRhs p f start = scanr op start
 where op x xs = f x 'union' if p x then xs else []

```

Finally, we will also need a function `scanlRhs` which does the same job as `scanrRhs` but in the opposite direction. The easiest way to define `scanlRhs` is in terms of `scanrRhs` and `reverse`.

```

scanlRhs p f start = reverse . scanrRhs p f start . reverse

```

We now return to the function `lookSet`.

```

lookSet :: Ord s =>
 (s -> Bool) ->
 (s -> [s]) ->
 (s -> [s]) ->
 (s, [s]) ->
 [s]
lookSet p f g (nt,rhs) = foldrRhs p f (g nt) rhs

```

The function `lookSet` makes use of `foldrRhs` to fold over a right-hand side. As stated above, the function `foldrRhs` continues processing a right-hand side only if it encounters a nonterminal symbol for which `p` (so `isEmpty` in the `lookSet` instance `lookaheadp`) holds. Thus, the set `g nt` (`follow?nt` in the `lookSet` instance `lookaheadp`) is only important for those right-hand sides for `nt` that consist of non-terminals that can all derive the empty sequence. We can now (assuming that the definitions of the auxiliary functions are given) use the function `lookaheadp` instance of `lookSet` to compute the lookahead sets of all productions.

```

look nt rhs = lookaheadp exGrammar (nt,rhs)

```

```

? look 'S' "AaS"
dba
? look 'S' "B"
dba
? look 'S' "CB"
d
? look 'A' "SC"

```

```

dba
? look 'A' ""
ad
? look 'B' "A"
dba
? look 'B' "b"
b
? look 'C' "D"
d
? look 'D' "d"
d

```

It is clear from this result that `exGrammar` is not an LL(1)-grammar. Let us have a closer look at how these lookahead sets are obtained. We will have to use the functions `firsts` and `follow` and the predicate `isEmpty` for computing intermediate results. The corresponding subsections explain how to compute these intermediate results.

For the lookahead set of the production  $A \rightarrow AaS$  we fold over the right-hand side  $AaS$ . Folding stops at 'a' and we obtain

```

firsts? 'A' 'union' firsts? 'a'
==
"dba" 'union' "a"
==
"dba"

```

For the lookahead set of the production  $A \rightarrow SC$  we fold over the right-hand side  $SC$ . Folding stops at  $C$  since it cannot derive the empty sequence, and we obtain

```

firsts? 'S' 'union' firsts? 'C'
==
"dba" 'union' "d"
==
"dba"

```

Finally, for the lookahead set of the production  $B \rightarrow A$  we fold over the right-hand side  $A$ . In this case we fold over the complete (one element) list and we obtain

```

firsts? 'A' 'union' follow? 'B'
==
"dba" 'union' "d"
==
"dba"

```

The other lookahead sets are computed in a similar way.

### 10.2.6. Implementation of empty

Many functions defined in this chapter make use of a predicate `isEmpty`, which tests whether or not the empty sequence can be derived from a nonterminal. This subsection defines this function. Consider the grammar `exGrammar`. We are now only interested in deriving sequences which contain only nonterminal symbols (since it is impossible to derive the empty string if a terminal occurs). Therefore we only have to consider the productions in which no terminal symbols appear in the right-hand sides.

$$\begin{aligned} S &\rightarrow B \mid CB \\ A &\rightarrow SC \mid \epsilon \\ B &\rightarrow A \\ C &\rightarrow D \end{aligned}$$

One can immediately see from those productions that the nonterminal  $A$  derives the empty string in one step. To know whether there are any nonterminals which derive the empty string in more than one step we eliminate the productions for  $A$  and we eliminate all occurrences of  $A$  in the right hand sides of the remaining productions

$$\begin{aligned} S &\rightarrow B \mid CB \\ B &\rightarrow \epsilon \\ C &\rightarrow D \end{aligned}$$

One can now conclude that the nonterminal  $B$  derives the empty string in two steps. Doing the same with  $B$  as we did with  $A$  gives us the following productions

$$\begin{aligned} S &\rightarrow \epsilon \mid C \\ C &\rightarrow D \end{aligned}$$

One can now conclude that the nonterminal  $S$  derives the empty string in three steps. Doing the same with  $S$  as we did with  $A$  and  $B$  gives us the following productions

$$C \rightarrow D$$

At this stage we can conclude that there are no more new nonterminals which derive the empty string.

We now give the Haskell implementation of the algorithm described above. The algorithm is iterative: it does the same steps over and over again until some desired condition is met. For this purpose we use function `fixedPoint`, which takes a function and a set, and repeatedly applies the function to the set, until the set does not change anymore.

```

fixedPoint :: Ord a => ([a] -> [a]) -> [a] -> [a]
fixedPoint f xs | xs == nexts = xs
 | otherwise = fixedPoint f nexts
 where nexts = f xs

```

`fixedPoint f` is sometimes called the *fixed-point* of `f`. Function `isEmpty` determines whether or not a nonterminal can derive the empty string. A nonterminal can derive the empty string if it is a member of the `emptySet` of a grammar.

```

isEmpty :: (Symbol s, Ord s) => CFG s -> s -> Bool
isEmpty grammar = ('elem' emptySet grammar)

```

The `emptySet` of a grammar is obtained by the iterative process described in the example above. We start with the empty set of nonterminals, and at each step  $n$  of the computation of the `emptySet` as a `fixedPoint`, we add the nonterminals that can derive the empty string in  $n$  steps. Function `emptyStepf` adds a nonterminal if there exists a production for the nonterminal of which all elements can derive the empty string.

```

emptySet :: (Symbol s, Ord s) => CFG s -> [s]
emptySet grammar = fixedPoint (emptyStepf grammar) []

emptyStepf :: (Symbol s, Ord s) => CFG s -> [s] -> [s]
emptyStepf grammar set =
 nub (map fst (filter (\(nt,rhs) -> all ('elem' set) rhs)
 (prods grammar)
)
)

```

### 10.2.7. Implementation of first and last

Function `firsts` takes a grammar, and returns for each symbol of the grammar (so also the terminal symbols) the set of terminal symbols with which a sentence derived from that symbol can start. The first set of a terminal symbol is the terminal symbol itself.

The set of firsts each symbol consists of that symbol itself, plus the (first) symbols that can be derived from that symbol in one or more steps. So the set of firsts can be computed by an iterative process, just as the function `isEmpty`.

Consider the grammar `exGrammar` again. We start the iteration with

```

[('S', "S"), ('A', "A"), ('B', "B"), ('C', "C"), ('D', "D")
, ('a', "a"), ('b', "b"), ('d', "d")
]

```

## 10. LL Parsing

Using the productions of the grammar we can derive in one step the following lists of first symbols.

```
[('S', "ABC"), ('A', "S"), ('B', "Ab"), ('C', "D"), ('D', "d")]
```

and the union of these two lists is

```
[('S', "SABC"), ('A', "AS"), ('B', "BAb"), ('C', "CD"), ('D', "Dd"),
, ('a', "a"), ('b', "b"), ('d', "d")]
```

In *two* steps we can derive

```
[('S', "SABd"), ('A', "ABC"), ('B', "S"), ('C', "d"), ('D', "")]
```

and again we have to take the union of this list with the previous result. We repeat this process until the list doesn't change anymore. For `exGrammar` this happens when:

```
[('S', "SABCDabd")
, ('A', "SABCDabd")
, ('B', "SABCDabd")
, ('C', "CDd")
, ('D', "Dd")
, ('a', "a")
, ('b', "b")
, ('d', "d")
]
```

Function `firsts` is defined as the `fixedPoint` of a step function that iterates the first computation one more step. The `fixedPoint` starts with the list that contains all symbols paired with themselves.

```
firsts :: (Symbol s, Ord s) => CFG s -> [(s,[s])]
firsts grammar =
 fixedPoint (firstStepf grammar) (startSingle grammar)

startSingle :: (Ord s, Symbol s) => CFG s -> [(s,[s])]
startSingle grammar = map (\x -> (x,[x])) (symbols grammar)
```

The step function takes the old approximation and performs one more iteration step. At each of these iteration steps we have to add the start list with which the iteration started again.

```

firstStepf :: (Ord s, Symbol s) =>
 CFG s -> [(s,[s])] -> [(s,[s])]
firstStepf grammar approx = (startSingle grammar)
 'combine' (compose (first1 grammar) approx)

combine :: Ord s => [(s,[s])] -> [(s,[s])] -> [(s,[s])]
combine xs = foldr insert xs
 where insert (a,bs) [] = [(a,bs)]
 insert (a,bs) ((c,ds):rest)
 | a == c = (a, union bs ds) : rest
 | otherwise = (c,ds) : (insert (a,bs) rest)

compose :: Ord a => [(a,[a])] -> [(a,[a])] -> [(a,[a])]
compose r1 r2 = [(a, unions (map (r2?) bs)) | (a,bs) <- r1]

```

Finally, function `first1` computes the direct first symbols of all productions, taking into account that some nonterminals can derive the empty string, and combines the results for the different nonterminals.

```

first1 :: (Symbol s, Ord s) => CFG s -> [(s,[s])]
first1 grammar =
 map (\(nt,fs) -> (nt,unions fs))
 (group (map (\(nt,rhs) -> (nt,foldrRhs (isEmpty grammar)
 single
 []
 rhs)
)
 (prods grammar)
)
)

```

where `group` groups elements with the same first element together

```

group :: Eq a => [(a,b)] -> [(a,[b])]
group = foldr insertPair []

insertPair :: Eq a => (a,b) -> [(a,[b])] -> [(a,[b])]
insertPair (a,b) [] = [(a,[b])]
insertPair (a,b) ((c,ds):rest) =
 if a==c then (c,(b:ds)):rest else (c,ds):(insertPair (a,b) rest)

```

function `single` takes an element and returns the set with the element, and `unions` returns the union of a set of sets.

Function `lasts` is defined using function `firsts`. Suppose we reverse the right-hand sides of all productions of a grammar. Then the set of firsts of this reversed grammar is the set of lasts of the original grammar. This idea is implemented in the following functions.

```
reverseGrammar :: Symbol s => CFG s -> CFG s
reverseGrammar =
 \ (s,al) -> (s,map (\(nt,rhs) -> (nt,reverse rhs)) al)

lasts :: (Symbol s, Ord s) => CFG s -> [(s,[s])]
lasts = firsts . reverseGrammar
```

### 10.2.8. Implementation of follow

The final function we have to implement is the function `follow`, which takes a grammar, and returns an association list in which nonterminals are associated to symbols that can follow upon the nonterminal in a derivation. A nonterminal `n` is associated to a list containing terminal `t` in `follow` if `n` and `t` follow each other in some sequence of symbols occurring in some leftmost derivation. We can compute pairs of such adjacent symbols by splitting up the right-hand sides with length at least 2 and, using `lasts` and `firsts`, compute the symbols which appear at the end resp. at the beginning of strings which can be derived from the left resp. right part of the split alternative. Our grammar `exGrammar` has three alternatives with length at least 2: "AaS", "CB" and "SC". Function `follow` uses the functions `firsts` and `lasts` and the predicate `isEmpty` for intermediate results. The previous subsections explain how to compute these functions.

Let's see what happens with the alternative "AaS". The lists of all nonterminal symbols that can appear at the end of sequences of symbols derived from "A" and "Aa" are "ADC" and "" respectively. The lists of all terminal symbols which can appear at the beginning of sequences of symbols derived from "aS" and "S" are "a" and "dba" respectively. Zipping together those lists shows that an 'A', a 'D' and a 'C' can be followed by an 'a'. Splitting the alternative "CB" in the middle produces sets of firsts and sets of lasts "CD" and "dba". Splitting the alternative "SC" in the middle produces sets of firsts and sets of lasts "SDCAB" and "d". From the first pair we can see that a 'C' and a 'D' can be followed by a 'd', a 'b' and an 'a'. From the second pair we see that an 'S', a 'D', a 'C', an 'A', and a 'B' can be followed by a 'd'. Combining all these results gives:

```
[('S',"d"),('A',"ad"),('B',"d"),('C',"adb"),('D',"adb")]
```

The function `follow` uses the functions `scanrAlt` and `scanlAlt`. The lists produced by these functions are exactly the ones we need: using the function `zip` from the



standard prelude we can combine the lists. For example: for the alternative "AaS" the functions `scanlAlt` and `scanrAlt` produce the following lists:

```
[[], "ADC", "DC", "SDCAB"]
["dba", "a", "dba", []]
```

Only the two middle elements of both lists are important (they correspond to the nontrivial splittings of "AaS"). Thus, we only have to consider alternatives of length at least 2. We start the computation of `follow` with assigning the empty follow set to each symbol:

```
follow :: (Symbol s, Ord s) => CFG s -> [(s,[s])]
follow grammar = combine (followNE grammar) (startEmpty grammar)

startEmpty grammar = map (\x -> (x,[])) (symbols grammar)
```

The real work is done by functions `followNE` and function `splitProds`. Function `followNE` passes the right arguments on to function `splitProds`, and removes all nonterminals from the set of firsts, and all terminals from the set of lasts. Function `splitProds` splits the productions of length at least 2, and pairs the last nonterminals with the first terminals.

```
followNE :: (Symbol s, Ord s) => CFG s -> [(s,[s])]
followNE grammar = splitProds
 (prods grammar)
 (isEmpty grammar)
 (isTfirsts grammar)
 (isNlasts grammar)
 where isTfirsts = map (\(x,xs) -> (x,filter isT xs)) . firsts
 isNlasts = map (\(x,xs) -> (x,filter isN xs)) . lasts

splitProds :: (Symbol s, Ord s) =>
 [(s,[s])] -> -- productions
 (s -> Bool) -> -- isEmpty
 [(s,[s])] -> -- terminal firsts
 [(s,[s])] -> -- nonterminal lasts
 [(s,[s])]

splitProds prods p fset lset =
 map (\(nt,rhs) -> (nt,nub rhs)) (group pairs)
 where pairs = [(l, f)
 | rhs <- map snd prods
 , length rhs >= 2
 , (fs, ls) <- zip (rightscan rhs) (leftscan rhs)]
```

## 10. LL Parsing

```
 , l <- ls
 , f <- fs
]
leftscan = scanLRhs p (lset?) []
rightscan = scanrRhs p (fset?) []
```

**Exercise 10.7.** Give the Rose tree representation of the parse tree corresponding to the derivation of the sentence `ccccba` using grammar `gramm1`.

**Exercise 10.8.** Give the Rose tree representation of the parse tree corresponding to the derivation of the sentence `abbb` using grammar `gramm2`' defined in the exercises of the previous section.

**Exercise 10.9.** Give the Rose tree representation of the parse tree corresponding to the derivation of the sentence `acbab` using grammar `gramm3`.

**Exercise 10.10.** In this exercise we will take a closer look at the functions `foldrRhs` and `scanrRhs` which are the essential ingredients of the implementation of the grammar analysis algorithms. From the definitions it is clear that grammar analysis is easily expressed via a calculus for (finite) sets. A calculus for finite sets is implicit in the programs for LL(1) parsing. Since the code in this module is obscured by several implementation details we will derive the functions `foldrRhs` and `scanrRhs` in a stepwise fashion. In this derivation we will use the following:

A (finite) set is implemented by a list with no duplicates. In order to construct a set, the following operations may be used:

|                     |                 |                              |                                           |
|---------------------|-----------------|------------------------------|-------------------------------------------|
| <code>[]</code>     | <code>::</code> | <code>[a]</code>             | the empty set of <code>a</code> -elements |
| <code>union</code>  | <code>::</code> | <code>[a] → [a] → [a]</code> | the union of two sets                     |
| <code>unions</code> | <code>::</code> | <code>[[a]] → [a]</code>     | the generalised union                     |
| <code>single</code> | <code>::</code> | <code>a → [a]</code>         | the singleton function                    |

These operations satisfy the well-known laws for set operations.

1. Define a function `list2Set :: [a] → [a]` which returns the set of elements occurring in the argument.
2. Define `list2Set` as a `foldr`.
3. Define a function `pref p :: [a] → [a]` which given a list `xs` returns the set of elements corresponding to the longest prefix of `xs` all of whose elements satisfy `p`.
4. Define a function `prefplus p :: [a] → [a]` which given a list `xs` returns the set of elements in the longest prefix of `xs` all of whose elements satisfy `p` together with the first element of `xs` that does not satisfy `p` (if this element exists at all).
5. Define `prefplus p` as a `foldr`.
6. Show that `prefplus p = foldrRhs p single []`.
7. It can be shown that

$$\text{foldrRhs } p \ f \ [] = \text{unions} \ . \ \text{map } f \ . \ \text{prefplus } p$$

for all set-valued functions `f`. Give an informal description of the functions `foldrRhs p f []` and `foldrRhs p f start`.

8. The standard function `scanr` is defined by

$$\text{scanr } f \text{ } q0 = \text{map } (\text{foldr } f \text{ } q0) \text{ } \cdot \text{ tails}$$

where `tails` is a function which takes a list `xs` and returns a list with all tailsegments (postfixes) of `xs` in decreasing length. The function `scanrRhs` is defined in a similar way

$$\text{scanrRhs } p \text{ } f \text{ } \text{start} = \text{map } (\text{foldrRhs } p \text{ } f \text{ } \text{start}) \text{ } \cdot \text{ tails}$$

Give an informal description of the function `scanrRhs`.

**Exercise 10.11.** The computation of the functions `empty` and `firsts` is not restricted to nonterminals only. For terminal symbols `s` these functions are defined by

$$\begin{aligned} \text{empty } s &= \text{False} \\ \text{firsts } s &= \{s\} \end{aligned}$$

Using the definitions in the previous exercise, compute the following.

1. For the example grammar `gramm1` and two of its productions  $A \rightarrow bSA$  and  $B \rightarrow Cb$ .
  - a) `foldrRhs empty firsts [] bSA`
  - b) `foldrRhs empty firsts [] Cb`
2. For the example grammar `gramm3` and its production  $S \rightarrow AaS$ 
  - a) `foldrRhs empty firsts [] AaS`
  - b) `scanrRhs empty firsts [] AaS`

## 10. LL Parsing

## 11. LL versus LR parsing

The parsers that were constructed using parser combinators in Chapter 3 are non-deterministic. They can recognize sentences described by ambiguous grammars by returning a *list of solutions*, rather than just one solution; each solution contains a parse tree and the part of the input that was left unprocessed. Nondeterministic parsers are of the following type:

```
type Parser a b = [a] -> [(b, [a])]
```

where `a` denotes the alphabet type, and `b` denotes the parse tree type.

In chapter 10 we turned to *deterministic* parsers. Here, parsers have only one result, consisting of a parse tree of type `b` and the remaining part of the input string:

```
type Parser a b = [a] -> (b, [a])
```

Ambiguous grammars are not allowed anymore. Also, in the case of parsing input containing syntax errors, we cannot return an empty list of successes anymore; instead there should be some mechanism of raising an error.

There are various algorithms for deterministic parsing. They impose some additional constraints on the form of the grammar: not every context free grammar is allowable by these algorithms. The parsing algorithms can be modelled by making use of a *stack machine*. There are two fundamentally different deterministic parsing algorithms:

- LL parsing, also known as *top-down parsing*
- LR parsing, also known as *bottom-up parsing*

The first ‘L’ in these acronyms stands for ‘Left-to-right’, that is, input is processed in the order it is read. So these algorithms are both suitable for reading an input stream, e.g. from a file. The second ‘L’ in ‘LL-parsing’ stands for ‘Leftmost derivation’, as the parsing mimics doing a leftmost derivation of a sentence (see Section 2.4). The ‘R’ in ‘LR-parsing’ stands for ‘Rightmost derivation’ of the sentences. The parsing algorithms are normally referred to as  $LL(k)$  or  $LR(k)$ , where  $k$  is the number of unread symbols that the parser is allowed to ‘look ahead’. In most practical cases,  $k$  is taken to be 1.

In Chapter 10 we studied the  $LL(1)$  parsing algorithm extensively, including the so-called  $LL(1)$ -property to which grammars must abide. In this chapter we start with an example application of the  $LL(1)$  algorithm. Next, we turn to the  $LR(1)$  algorithm.

## 11.1. LL(1) parser example

The  $LL(1)$  parsing algorithm was implemented in Section 10.2.3. Function `ll1` defined there takes a grammar and an input string, and returns a single result consisting of a parse tree and rest string. The function was implemented by calling a generalized function named `gll1`; generalized in the sense that the function takes an additional list of symbols that need to be recognized. That generalization is then called with a singleton list containing only the root nonterminal.

### 11.1.1. An LL(1) checker

Here, we define a slightly simplified version of the algorithm: it doesn't return a parse tree, so it need not be concerned with building it; it merely *checks* whether or not the input string is a sentence. Hence the result of the function is simply a boolean value:

```
check :: String -> Bool
check input = run ['S'] input
```

As in chapter 10, the function is implemented by calling a generalized function `run` with a singleton containing just the root nonterminal. Now the function `run` takes, in addition to the input string, a list of symbols (which is initially called with the above-mentioned singleton). That list is used as some kind of stack, as elements are prepended at its front, and removed at the front in other occasions. Therefore we refer to the whole operation as a *stack machine*, and that's why the function is named `run`: it *runs* the stack machine. Function `run` is defined as follows:

```
run :: Stack -> String -> Bool
run [] [] = True
run [] (x:xs) = False
run (a:as) input | isT a = not(null input)
 && a==hd input
 && run as (tl input)
 | isN a = run (rs++as) input
 where rs = select a (hd input)
```

So, when called with an empty stack and an empty string, the function succeeds. If the stack is empty, but the input is not, it fails, as there is junk input present. If the stack is nonempty, case distinction is done on `a`, the top of the stack. If it is a terminal, the input should begin with it, and the machine is called recursively for the rest of the input. In the case of a nonterminal, we push `rs` on the stack, and leave the input unchanged in the recursive call. This is a simple tail recursive function, and could imperatively be implemented in a single loop that runs while the stack is non-empty.

The new symbols `rs` that are pushed on the stack is the right hand side of a production rule for nonterminal `a`. It is selected from the available alternatives by function

**select**. For making the choice, the first input symbol is passed to **select** as well. This is where you can see that the algorithm is *LL(1)*: it is allowed to look ahead *one* symbol.

```
This is select the alternative (since it is selected filter ok . prods) gram
 where ok p@(n,_) = n==a && x 'elem' lahP gram p
```

So, from all productions of the grammar returned by **prods**, the **ok** ones are taken, of which there should be only one (this is ensured by the *LL(1)*-property); that single production is retrieved by **hd**, and of it only the right hand side is needed, hence the call to **snd**. Now a production is **ok**, if the nonterminal **n** is **a**, the one we are looking for, and moreover the first input symbol **x** is member of the *lookahead set* of this production.

Determining the lookahead sets of all productions of a grammar is the tricky part of the *LL(1)* parsing algorithm. It is described in Sections 10.2.5 to 10.2.8. Though the general formulation of the algorithm is quite complex, its application in a simple case is rather intuitive. So let's study an example grammar: arithmetic expressions with operators of two levels of precedence.

### 11.1.2. An *LL(1)* grammar

The idea for the grammar for arithmetical expressions was given in Section 2.5.7: we need auxiliary notions of 'Term' and 'Factor' (actually, we need as much additional notions as there are levels of precedence). The most straightforward definition of the grammar is shown in the left column below.

Unfortunately, this grammar doesn't abide to the *LL(1)*-property. The reason for this is that the grammar contains rules with common prefixes on the right hand side (for *E* and *T*). A way out is the application of a grammar transformation known as *left factoring*, as described in Section 2.5.4. The result is a grammar where there is only one rule for *E* and *T*, and the non-common part of the right hand sides is described by additional nonterminals, *P* and *M*. The resulting grammar is shown in the right column below. For being able to treat end-of-input as if it were a character, we also add an additional rule, which says that the input consists of an expression followed by end-of-input (designated as '#' here).

## 11. LL versus LR parsing

|                       |                          |
|-----------------------|--------------------------|
| $E \rightarrow T$     | $S \rightarrow E \#$     |
| $E \rightarrow T + E$ | $E \rightarrow TP$       |
| $T \rightarrow F$     | $P \rightarrow \epsilon$ |
| $T \rightarrow F * T$ | $P \rightarrow + E$      |
| $F \rightarrow N$     | $T \rightarrow FM$       |
| $F \rightarrow ( E )$ | $M \rightarrow \epsilon$ |
| $N \rightarrow 1$     | $M \rightarrow * T$      |
| $N \rightarrow 2$     | $F \rightarrow N$        |
| $N \rightarrow 3$     | $F \rightarrow ( E )$    |
|                       | $N \rightarrow 1$        |
|                       | $N \rightarrow 2$        |
|                       | $N \rightarrow 3$        |

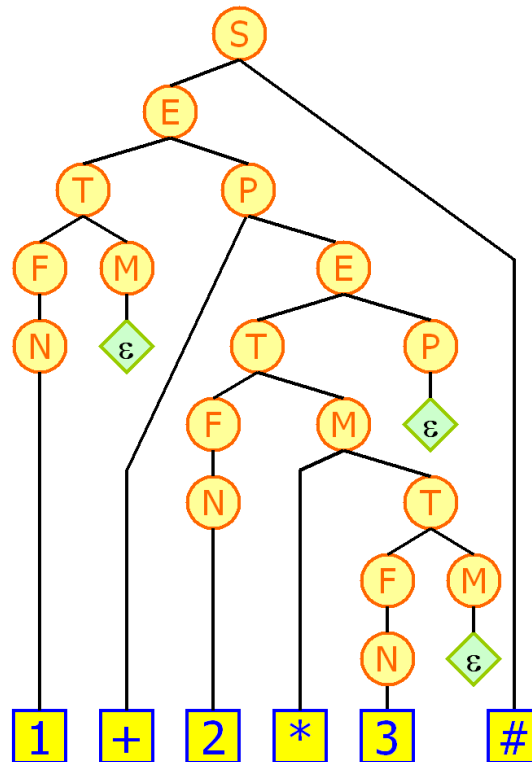
For determining the lookahead sets of each nonterminal, we also need to analyze whether a nonterminal can produce the empty string, and which terminals can be the first symbol of any string produced by each nonterminal. These properties are named *empty* and *first* respectively. All properties for the example grammar are summarized in the table below. Note that *empty* and *first* are properties of a nonterminal, whereas *lookahead* is a property of a single production rule.

| production               | <i>empty</i> | <i>first</i> | <i>lookahead</i> |
|--------------------------|--------------|--------------|------------------|
| $S \rightarrow E \#$     | no           | ( 1 2 3      | first(E) ( 1 2 3 |
| $E \rightarrow TP$       | no           | ( 1 2 3      | first(T) ( 1 2 3 |
| $P \rightarrow \epsilon$ | yes          | +            | follow(P) ) #    |
| $P \rightarrow + E$      |              |              | immediate +      |
| $T \rightarrow FM$       | no           | ( 1 2 3      | first(F) ( 1 2 3 |
| $M \rightarrow \epsilon$ | yes          | *            | follow(M) + ) #  |
| $M \rightarrow * T$      |              |              | immediate *      |
| $F \rightarrow N$        | no           | ( 1 2 3      | first(N) 1 2 3   |
| $F \rightarrow ( E )$    |              |              | immediate (      |
| $N \rightarrow 1$        |              |              | immediate 1      |
| $N \rightarrow 2$        |              |              | immediate 2      |
| $N \rightarrow 3$        |              |              | immediate 3      |

### 11.1.3. Using the LL(1) parser

A sentence of the language described by the example grammar is  $1+2*3$ . Because of the high precedence of  $*$  relative to  $+$ , the parse tree should reflect that the 2 and 3 should be multiplied, rather than the  $1+2$  and 3. Indeed, the parse tree does so:





Now when we do a step-by-step analysis of how the parse tree is constructed by the stack machine, we notice that the parse tree is traversed in a depth-first fashion, where the left subtrees are analysed first. Each node corresponds to the application of a production rule. The order in which the production rules are applied is a pre-order traversal of the tree. The tree is constructed top-down: first the root is visited, and then each time the first remaining nonterminal is expanded. From the table in Figure 11.1 it is clear that the contents of the stack describes what is still to be expected on the input. Initially, of course, this is the root nonterminal  $S$ . Whenever a terminal is on top of the stack, the corresponding symbol is read from the input.

## 11.2. LR parsing

Another algorithm for doing deterministic parsing using a stack machine, is  $LR(1)$ -parsing. Actually, what is described in this section is known as *Simple LR* parsing or  $SLR(1)$ -parsing. It is still rather complicated, though.

A nice property of LR-parsing is that it is in many ways exactly the opposite, or *dual*, of LL-parsing. Some of these ways are:

- LL does a leftmost derivation, LR does a rightmost derivation

## 11. LL versus LR parsing

- LL *starts* with only the root nonterminal on the stack, LR *ends* with only the root nonterminal on the stack
- LL *ends* when the stack is empty, LR *starts* with an empty stack
- LL uses the stack for designating what is still to be *expected*, LR uses the stack for designating what has already been *seen*
- LL builds the parse tree *top down*, LR builds the parse tree *bottom up*
- LL continuously pops a nonterminal off the stack, and pushes a corresponding right hand side; LR tries to recognize a right hand side on the stack, pops it, and pushes the corresponding nonterminal
- LL thus *expands* nonterminals, while LR *reduces* them
- LL reads terminal when it pops one *off* the stack, LR reads terminals while it pushes them *on* the stack
- LL uses grammar rules in an order which corresponds to *pre-order* traversal of the parse tree, LR does a *post-order* traversal.

### 11.2.1. A stack machine for *SLR* parsing

As in Section 10.1.1 a stack machine is used to parse sentences. A stack machine for a grammar  $G$  has a stack and an input. When the parsing starts the stack is empty. The stack machine performs one of the following two actions.

1. **Shift:** Move the first symbol of the remaining input to the top of the stack.
2. **Reduce:** Choose a production rule  $N \rightarrow \alpha$ ; pop the sequence  $\alpha$  from the top of the stack; push  $N$  onto the stack.

These actions are performed until the stack only contains the start symbol. A stack machine for  $G$  accepts an input if it can terminate with only the start symbol on the stack when the whole input has been processed. Let us take the example of Section 10.1.1.

| stack | input |
|-------|-------|
|       | aab   |
| a     | ab    |
| aa    | b     |
| baa   |       |
| Saa   |       |
| Sa    |       |
| S     |       |

Note that with our notation the top of the stack is at the left side. To decide which production rule is to be used for the reduction, the symbols on the stack must be read in reverse order.

The stack machine for  $G$  accepts the input string `aab` because it terminates with the start symbol on the stack and the input is completely processed. The stack machine performs three shift actions and then three reduce actions.

The SLR parser only performs a reduction with a grammar rule  $N \rightarrow \alpha$  if the next symbol on the input is in the follow set of  $N$ .

**Exercise 11.1.** Give for `gramm1` of Section 10.1.2 all the states of the stack machine for a derivation of the input `ccccba`.

**Exercise 11.2.** Give for `gramm2` of Section 10.1.2 all the states of the stack machine for a derivation of the input `abbb`.

**Exercise 11.3.** Give for `gramm3` of Section 10.1.2 all the states of the stack machine for a derivation of the input `acbab`.

## 11.3. LR parse example

### 11.3.1. An LR checker

for a start, the main function for LR parsing calls a generalized stack function with an empty stack:

```
check' :: String -> Bool
check' input = run' [] input
```

The stack machine terminates when it finds the stack containing just the root non-terminal. Otherwise it either pushes the first input symbol on the stack ('Shift'), or it drops some symbols off the stack (which should be the right hand side of a rule) and pushes the corresponding nonterminal ('Reduce').

```
run' :: Stack -> String -> Bool
run' ['S'] [] = True
run' ['S'] (x:xs) = False
run' stack (x:xs) = case action of
 Shift -> run' (x: stack) xs
 Reduce a n -> run' (a:drop n stack) (x:xs)
 Error -> False
 where action = select' stack x
```

In the case of LL-parsing the hard part was selecting the right rule to expand; here we have the hard decision of whether to reduce according to a (and which?) rule, or to shift the next symbol. This is done by the `select'` function, which is allowed to inspect the first input symbol `x` and the *entire* stack: after all, it needs to find the right hand side of a rule on the stack.

In the right column in Figure 11.1 the LR derivation of sentence `1+2*3` is shown. Compare closely to the left column, which shows the LL derivation, and note the duality of the processes.

11. LL versus LR parsing

| LL derivation            |       |        |           | LR derivation            |       |          |           |
|--------------------------|-------|--------|-----------|--------------------------|-------|----------|-----------|
| rule                     | read  | ↓stack | remaining | rule                     | read  | stack↓   | remaining |
|                          |       | $S$    | $1+2*3$   |                          |       |          | $1+2*3$   |
| $S \rightarrow E$        |       | $E$    | $1+2*3$   | shift                    | 1     | 1        | $+2*3$    |
| $E \rightarrow TP$       |       | $TP$   | $1+2*3$   | $N \rightarrow 1$        | 1     | $N$      | $+2*3$    |
| $T \rightarrow FM$       |       | $FMP$  | $1+2*3$   | $F \rightarrow N$        | 1     | $F$      | $+2*3$    |
| $F \rightarrow N$        |       | $NMP$  | $1+2*3$   | $M \rightarrow \epsilon$ | 1     | $FM$     | $+2*3$    |
| $N \rightarrow 1$        |       | $1MP$  | $1+2*3$   | $T \rightarrow FM$       | 1     | $T$      | $+2*3$    |
| read                     | 1     | $MP$   | $+2*3$    | shift                    | 1+    | $T+$     | $2*3$     |
| $M \rightarrow \epsilon$ | 1     | $P$    | $+2*3$    | shift                    | 1+2   | $T+2$    | $*3$      |
| $P \rightarrow +E$       | 1     | $+E$   | $+2*3$    | $N \rightarrow 2$        | 1+2   | $T+N$    | $*3$      |
| read                     | 1+    | $E$    | $2*3$     | $F \rightarrow N$        | 1+2   | $T+F$    | $*3$      |
| $E \rightarrow TP$       | 1+    | $TP$   | $2*3$     | shift                    | 1+2*  | $T+F*$   | 3         |
| $T \rightarrow FM$       | 1+    | $FMP$  | $2*3$     | shift                    | 1+2*3 | $T+F*3$  |           |
| $F \rightarrow N$        | 1+    | $NMP$  | $2*3$     | $N \rightarrow 3$        | 1+2*3 | $T+F*N$  |           |
| $N \rightarrow 2$        | 1+    | $2MP$  | $2*3$     | $F \rightarrow N$        | 1+2*3 | $T+F*F$  |           |
| read                     | 1+2   | $MP$   | $*3$      | $M \rightarrow \epsilon$ | 1+2*3 | $T+F*FM$ |           |
| $M \rightarrow *T$       | 1+2   | $*TP$  | $*3$      | $T \rightarrow FM$       | 1+2*3 | $T+F*T$  |           |
| read                     | 1+2*  | $TP$   | 3         | $M \rightarrow *T$       | 1+2*3 | $T+FM$   |           |
| $T \rightarrow FM$       | 1+2*  | $FMP$  | 3         | $T \rightarrow FM$       | 1+2*3 | $T+T$    |           |
| $F \rightarrow N$        | 1+2*  | $NMP$  | 3         | $P \rightarrow \epsilon$ | 1+2*3 | $T+TP$   |           |
| $N \rightarrow 3$        | 1+2*  | $3MP$  | 3         | $E \rightarrow TP$       | 1+2*3 | $T+E$    |           |
| read                     | 1+2*3 | $MP$   |           | $P \rightarrow +E$       | 1+2*3 | $TP$     |           |
| $M \rightarrow \epsilon$ | 1+2*3 | $P$    |           | $E \rightarrow TP$       | 1+2*3 | $E$      |           |
| $P \rightarrow \epsilon$ | 1+2*3 |        |           | $S \rightarrow E$        | 1+2*3 | $S$      |           |

Figure 11.1.: LL and LR derivations of  $1+2*3$

### 11.3.2. LR action selection

The choice whether to shift or to reduce is made by function `select'`. It is defined as follows:

```
select' as x
| null items = Error
| null redItems = Shift
| otherwise = Reduce a (length rs)
 where items = dfa as
 redItems = filter red items
 (a,rs,_) = hd redItems
```

In the selection process a set (list) of so-called *items* plays a role. If the set is empty, there is an error. If the set contains at least one item, we filter the `red`, or *reducible* items from it. There should be only one, if the grammar has the LR-property. (Or rather: it is the LR-property that there is only one element in this situation). The reducible item is the production rule that can be reduced by the stack machine.

Now what are these items? An *item* is defined to be a production rule, augmented with a 'cursor position' somewhere in the right hand side of the rule. So for the rule  $F \rightarrow (E)$ , there are four possible items:  $F \rightarrow \cdot(E)$ ,  $F \rightarrow (\cdot E)$ ,  $F \rightarrow (E \cdot)$  and  $F \rightarrow (E)\cdot$ , where  $\cdot$  denotes the cursor.

for the rule  $F \rightarrow 2$  we have two items: one having the cursor in front of the single symbol, and one with the cursor after it:  $F \rightarrow \cdot 2$  and  $F \rightarrow 2 \cdot$ . For epsilon-rules there is a single item, where the cursor is at position 0 in the right hand side.

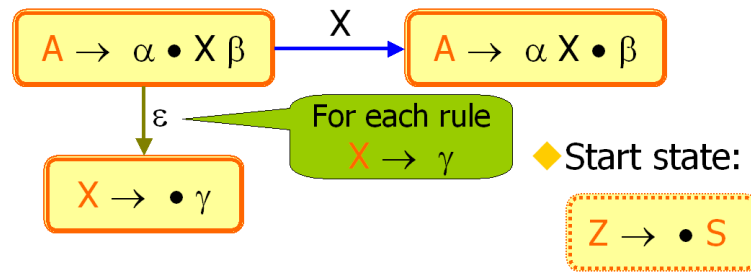
In the Haskell representation, the cursor is an integer which is added as a third element of the tuple, which already contains nonterminal and right hand side of the rule.

The items thus constructed are taken to be the states of a NFA (Nondeterministic Finite-state Automaton), as described in Section 8.1.2. We have the following transition relations in this NFA:

- The cursor can be 'stepped' to the next position. This transition is labeled with the symbol that is hopped over
- If the cursor is on front of a nonterminal, we can jump to an item describing the application of that nonterminal, where the cursor is at the beginning. This relation is an epsilon-transition of the NFA.

11. LL versus LR parsing

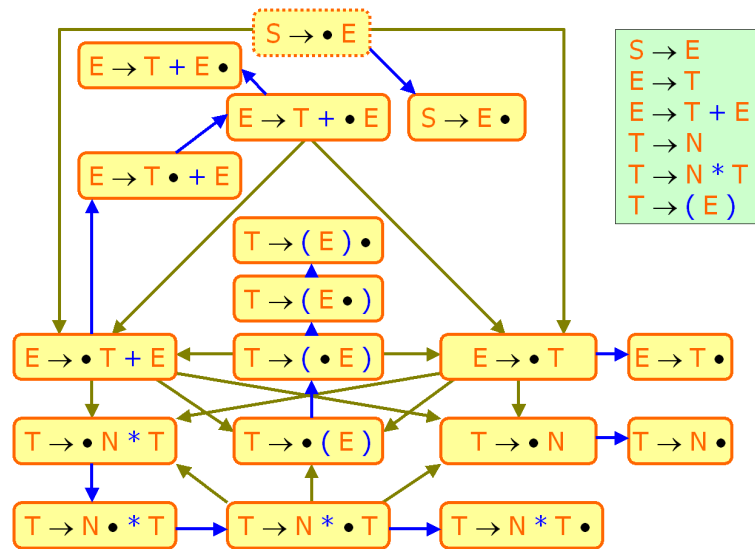
◆ NFA transitions:



As an example, let's consider an simplification of the arithmetic expression grammar:

- $E \rightarrow T$
- $E \rightarrow T + E$
- $T \rightarrow N$
- $T \rightarrow N * T$
- $T \rightarrow ( E )$

it skips the 'factor' notion as compared to the grammar earlier in this chapter, so it misses some well-formed expressions, but it serves only as an example for creating the states here. There are 18 states in this machine, as depicted here:

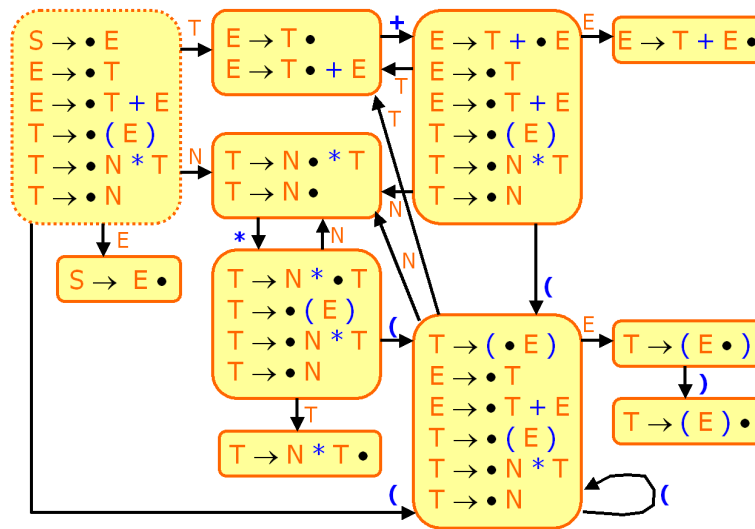


**Exercise 11.4** (no answer provided). How can you predict the number of states from inspecting the grammar?

**Exercise 11.5** (no answer provided). In the picture, the transition arrow labels are not shown. Add them.

**Exercise 11.6** (no answer provided). What makes this FA nondeterministic?

As was shown in Section 8.1.4, another automaton can be constructed that is deterministic (a DFA). That construction involves defining states which are *sets* of the original states. So in this case, the states of the DFA are *sets of items* (where items are rules with a cursor position). In the worst case we would have  $2^{18}$  states for the DFA-version of our example NFA, but it turns out that we are in luck: there are only 11 states in the DFA. Its states are rather complicated to depict, as they are sets of items, but it can be done:



**Exercise 11.7** (no answer provided). Check that this FA is indeed deterministic.

Given this DFA, let's return to our function that selects whether to shift or to reduce:

```
select' as x
| null items = Error
| null redItems = Shift
| otherwise = Reduce a (length rs)
 where items = dfa as
 redItems = filter red items
 (a,rs,_) = hd redItems
```

It runs the DFA, *using the contents of the stack (as) as input*. From the state where the DFA ends, which is by construction a set of items, the 'red' ones are filtered out. An item is 'red' (that is: reducible) if the cursor is at the end of the rule.

## 11. LL versus LR parsing

**Exercise 11.8** (no answer provided). Which of the states in the picture of the DFA contain red items? How many?

This is not the only condition that makes an item reducible; the second condition is that the lookahead input symbol is in the *follow* set of the nonterminal that is reduced to. Function `follow` was also needed in the LL analysis in Section 10.2.8. Both conditions are in the formal definition:

```
... where red (a,r,c) = c==length r && x 'elem' follow a
```

**Exercise 11.9** (no answer provided). How does this condition reflect that ‘the cursor is at the end’?

### 11.3.3. LR optimizations and generalizations

The stack machine `run`, by its nature, pushes and pops the stack continuously, and does a recursive call afterwards. In each call, for making the shift/reduce decision, the DFA is run on the (new) stack. In practice, it is considered a waste of time to do the full DFA transitions each time, as most of the stack remains the same after some popping and pushing at the top. Therefore, as an optimization, at each stack position, the corresponding DFA state is also stored. The states of the DFA can easily be numbered, so this amounts to just storing extra integers on the stack, tupled with the symbols that used to be on the stack. (An artificial bottom element should be placed on the stack initially, containing a dummy symbol and the number of the initial DFA state).

By analysing the grammar, two tables can be precomputed:

- Shift, that decides what is the new state when pushing a terminal symbol on the stack. This basically is the transition relation on the DFA.
- Action, that decides what action to take from a given state seeing a given input symbol.

Both tables can be implemented as a two-dimensional table of integers, of dimensions the number of states (typically, 100s to 1000s) times the number of symbols (typically, under 100).

As said in the introduction, the algorithm described here is a mere *Simple* LR parsing, or SLR(1). Its simplicity is in the reducibility test, which says:

```
... where red (a,r,c) = c==length r && x 'elem' follow a
```

The *follow* set is a rough approximation of what might follow a given nonterminal. But this set is not dependent of the context in which the nonterminal is used; maybe, in some contexts, the set is smaller. So, the SLR `red` function, may designate items as reducible, where it actually should not. For some grammars this might lead to a decision to reduce, where it should have done a shift, with a failing parser as a consequence.



An improvement, leading to a wider class of grammars that are allowable, would be to make the *follow* set context dependent. This means that it should vary for each *item* instead of for each *nonterminal*. It leads to a dramatic increase of the number of states in the DFA. And the full power of LR parsing is rarely needed.

A compromise position is taken by the so-called *LALR* parsers, or *Look Ahead LR* parsing. (A rather silly name, as *all* parsers look ahead...). In LALR parsers, the follow sets are context dependent, but when states in the DFA differ only with respect to the follow-part of their set-members (and not with respect to the item-part of them), the states are merged. Grammars that do not give rise to shift/reduce conflicts in this situation are said to be LALR-grammars. It is not really a natural notion; rather, it is a performance hack that gets most of LR power while keeping the size of the goto- and action-tables reasonable.

A widely used parser generator tool named *yacc* (for ‘yet another compiler compiler’) is based on an LALR engine. It comes with Unix, and was originally created for implementing the first C compilers. A commonly used clone of *yacc* is named *Bison*. *Yacc* is designed for doing the context-free aspects of analysing a language. The micro structure of identifiers, numbers, keywords, comments etc. is not handled by this grammar. Instead, it is described by regular expressions, which are analysed by a accompanying tool to *yacc* named *lex* (for ‘lexical scanner’). *Lex* is a preprocessor to *yacc*, that subdivides the input character stream into a stream of meaningful *tokens*, such as numbers, identifiers, operators, keywords etc.

### Bibliographical notes

The example DFA and NFA for LR parsing, and part of the description of the algorithm were taken from course notes by Alex Aiken and George Necula, which can be found at [www.cs.wright.edu/~tkprasad/courses/cs780](http://www.cs.wright.edu/~tkprasad/courses/cs780)

## 11. LL versus LR parsing

# Bibliography

- [1] A.V. Aho, Sethi R., and J.D. Ullman. *Compilers — Principles, Techniques and Tools*. Addison-Wesley, 1986.
- [2] R.S. Bird. Using circular programs to eliminate multiple traversals of data. *Acta Informatica*, 21:239–250, 1984.
- [3] R.S. Bird and P. Wadler. *Introduction to Functional Programming*. Prentice Hall, 1988.
- [4] W.H. Burge. Parsing. In *Recursive Programming Techniques*. Addison-Wesley, 1975.
- [5] J. Fokker. Functional parsers. In J. Jeuring and E. Meijer, editors, *Advanced Functional Programming*, volume 925 of *Lecture Notes in Computer Science*. Springer-Verlag, 1995.
- [6] R. Harper. Proof-directed debugging. *Journal of Functional Programming*, 1999. To appear.
- [7] G. Hutton. Higher-order functions for parsing. *Journal of Functional Programming*, 2(3):323 – 343, 1992.
- [8] B.W. Kernighan and R. Pike. Regular expressions — languages, algorithms, and software. *Dr. Dobb's Journal*, April:19 – 22, 1999.
- [9] D.E. Knuth. Semantics of context-free languages. *Math. Syst. Theory*, 2(2):127–145, 1968.
- [10] Niklas Røjemo. *Garbage collection and memory efficiency in lazy functional languages*. PhD thesis, Chalmers University of Technology, 1995.
- [11] S. Sippu and E. Soisalon-Soininen. *Parsing Theory, Vol. 1: Languages and Parsing*, volume 15 of *EATCS Monographs on THEoretical Computer Science*. Springer-Verlag, 1988.
- [12] S.D. Swierstra and P.R. Azero Alcocer. Fast, error correcting parser combinators: a short tutorial. In *SOFSEM'99*, 1999.
- [13] P. Wadler. How to replace failure by a list of successes: a method for exception handling, backtracking, and pattern matching in lazy functional languages. In J.P. Jouannaud, editor, *Functional Programming Languages and Computer Architecture*, pages 113 – 128. Springer, 1985. LNCS 201.

## Bibliography

## A. The Stack module

```
module Stack
 (Stack
 , emptyStack
 , isEmptyStack
 , push
 , pushList
 , pop
 , popList
 , top
 , split
 , mystack
)

where

data Stack x = MkS [x] deriving (Show,Eq)

emptyStack :: Stack x
emptyStack = MkS []

isEmptyStack :: Stack x -> Bool
isEmptyStack (MkS xs) = null xs

push :: x -> Stack x -> Stack x
push x (MkS xs) = MkS (x:xs)

pushList :: [x] -> Stack x -> Stack x
pushList xs (MkS ys) = MkS (xs ++ ys)

pop :: Stack x -> Stack x
pop (MkS xs) = if isEmptyStack (MkS xs)
 then error "pop on emptyStack"
 else MkS (tail xs)

popIf :: Eq x => x -> Stack x -> Stack x
```

## A. The Stack module

```
popIf x stack = if top stack == x
 then pop stack
 else error "argument and top of stack don't match"
```

```
popList :: Eq x => [x] -> Stack x -> Stack x
popList xs stack = foldr popIf stack (reverse xs)
```

```
top :: Stack x -> x
top (MkS xs) = if isEmptyStack (MkS xs)
 then error "top on emptyStack"
 else head xs
```

```
split :: Int -> Stack x -> ([x], Stack x)
split 0 stack = ([], stack)
split n (MkS []) = error "attempt to split the emptyStack"
split (n+1) (MkS (x:xs)) = (x:ys, stack')
 where
 (ys, stack') = split n (MkS xs)
```

```
mystack = MkS [1,2,3,4,5,6,7]
```

## B. Answers to exercises

**2.1** Three of the four strings are elements of  $L^*$ : abaabaaabaa, aaaabaaaa, baaaaabaa.

**2.2**  $\{\varepsilon\}$ .

**2.3**

$$\begin{aligned} & \emptyset L \\ = & \{ \text{Definition of concatenation of languages} \} \\ & \{ st \mid s \in \emptyset, t \in L \} \\ = & \{ s \in \emptyset \} \\ & \emptyset \end{aligned}$$

The other equalities can be proved in a similar fashion.

**2.4** The star operator on sets injects the elements of a set in a list; the star operator on languages concatenates the sentences of the language. The former star operator preserves more structure.

**2.5** Section 2.1 contains an inductive definition of the set of sequences over an arbitrary set  $X$ . Syntactical definitions for such sets follow immediately from this.

1. A grammar for  $X = \{a\}$  is given by

$$\begin{aligned} S & \rightarrow \varepsilon \\ S & \rightarrow aS \end{aligned}$$

2. A grammar for  $X = \{a, b\}$  is given by

$$\begin{aligned} S & \rightarrow \varepsilon \\ S & \rightarrow XS \\ X & \rightarrow a \mid b \end{aligned}$$

**2.6** A context free grammar for  $L$  is given by

$$\begin{aligned} S & \rightarrow \varepsilon \\ S & \rightarrow aSb \end{aligned}$$

**2.7** Analogous to the construction of the grammar for PAL.

$$\begin{array}{l}
 P \rightarrow \varepsilon \\
 | \text{ a} \\
 | \text{ b} \\
 | \text{ aPa} \\
 | \text{ bPb}
 \end{array}$$

**2.8** Analogous to the construction of PAL.

$$\begin{array}{l}
 M \rightarrow \varepsilon \\
 | \text{ aMa} \\
 | \text{ bMb}
 \end{array}$$

**2.9** First establish an inductive definition for parity sequences. An example of a grammar that can be derived from the inductive definition is:

$$P \rightarrow \varepsilon \mid 1P1 \mid 0P \mid P0$$

There are many other solutions.

**2.10** Again, establish an inductive definition for  $L$ . An example of a grammar that can be derived from the inductive definition is:

$$S \rightarrow \varepsilon \mid \text{aSb} \mid \text{bSa} \mid SS$$

Again, there are many other solutions.

**2.11** A sentence is a sentential form consisting only of terminals which can be derived in *zero* or more derivation steps from the start symbol (to be more precise: the sentential form consisting only of the start symbol). The start symbol is a non-terminal. The nonterminals of a grammar do not belong to the alphabet (the set of terminals) of the language we describe using the grammar. Therefore the start symbol cannot be a sentence of the language. As a consequence we have to perform at least *one* derivation step from the start symbol before we end up with a sentence of the language.

**2.12** The language consisting of the empty string only, i.e.,  $\{\varepsilon\}$ .

**2.13** This grammar generates the empty language, i.e.,  $\emptyset$ . In general, grammars such as this one that have no production rules without nonterminals on the right hand side, cannot produce any sentences with only terminal symbols. Each derivation will always contain nonterminals, so no sentences can be derived.

**2.14** The sentences in this language consist of zero or more concatenations of **ab**, i.e., the language is the set  $\{\mathbf{ab}\}^*$ .



**2.15** Yes. Each finite language is context free. A context free grammar can be obtained by taking one nonterminal and adding a production rule for each sentence in the language. For the language in Exercise 2.1, this procedure yields

$$\begin{aligned} S &\rightarrow ab \\ S &\rightarrow aa \\ S &\rightarrow baa \end{aligned}$$

**2.16** To bring the grammar into the form where we can directly apply the rule for associative separators, we introduce a new nonterminal:

$$\begin{aligned} A &\rightarrow AaA \\ A &\rightarrow B \\ B &\rightarrow b \mid c \end{aligned}$$

Now we can remove the ambiguity:

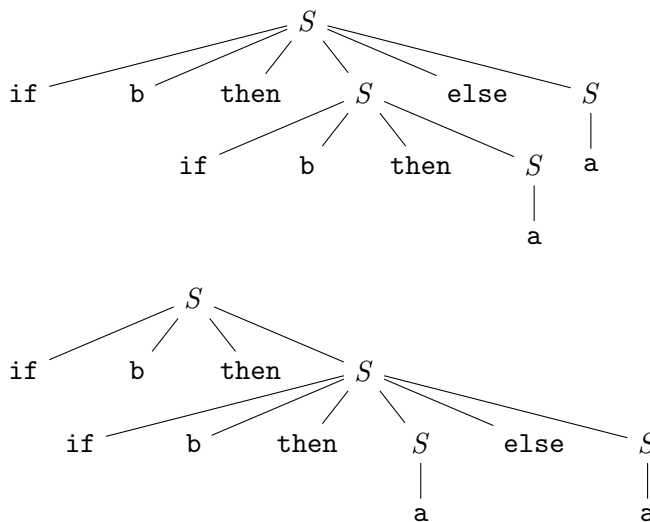
$$\begin{aligned} A &\rightarrow BaA \\ A &\rightarrow B \\ B &\rightarrow b \mid c \end{aligned}$$

It is now (optionally) possible to undo the auxiliary step of introducing the additional nonterminal by applying the rules for substituting right hand sides for nonterminal and removing unreachable productions. We then obtain:

$$\begin{aligned} A &\rightarrow baA \mid caA \\ A &\rightarrow b \mid c \end{aligned}$$

**2.17**

- Here are two parse trees for the sentence `if b then if b then a else a`:



B. Answers to exercises

- The rule we apply is: match **else** with the closest previous unmatched **if**. This means we prefer the second of the two parse trees above. The disambiguating rule is incorporated directly into the grammar:

$$\begin{array}{l}
 S \quad \quad \quad \rightarrow \textit{MatchedS} \mid \textit{UnmatchedS} \\
 \textit{MatchedS} \quad \rightarrow \text{if } \mathbf{b} \text{ then } \textit{MatchedS} \text{ else } \textit{MatchedS} \\
 \quad \quad \quad \quad \quad \mid \mathbf{a} \\
 \textit{UnmatchedS} \rightarrow \text{if } \mathbf{b} \text{ then } S \\
 \quad \quad \quad \quad \quad \mid \text{if } \mathbf{b} \text{ then } \textit{MatchedS} \text{ else } \textit{UnmatchedS}
 \end{array}$$

- An **else** clause is always matched with the closest previous unmatched **if**.

**2.18** An equivalent grammar for bit lists is

$$\begin{array}{l}
 L \rightarrow B Z \mid B \\
 Z \rightarrow , L Z \mid , L \\
 B \rightarrow 0 \mid 1
 \end{array}$$

**2.19**

- The grammar generates the language  $\{a^{2^n}b^m \mid m, n \in \mathbb{N}\}$ .
- An equivalent non left recursive grammar is

$$\begin{array}{l}
 S \rightarrow AB \\
 A \rightarrow \varepsilon \mid \mathbf{aa}A \\
 B \rightarrow \varepsilon \mid \mathbf{b}B
 \end{array}$$

**2.20** Of course, we can choose how we want to represent the different operators in concrete syntax. Choosing standard symbols, this is one possibility:

$$\begin{array}{l}
 \textit{Expr} \rightarrow \textit{Expr} + \textit{Expr} \\
 \quad \quad \mid \textit{Expr} * \textit{Expr} \\
 \quad \quad \mid \textit{Int}
 \end{array}$$

where *Int* is a nonterminal that produces integers. Note that the grammar given above is ambiguous. We could also give an unambiguous version, for example by introducing operator priorities.

**2.21** Recall the grammar for palindromes from Exercise 2.7, now with names for the productions:

$$\begin{array}{l}
 \text{Empty: } P \rightarrow \varepsilon \\
 \text{A: } \quad P \rightarrow \mathbf{a} \\
 \text{B: } \quad P \rightarrow \mathbf{b} \\
 \text{A}_2: \quad P \rightarrow \mathbf{aPa} \\
 \text{B}_2: \quad P \rightarrow \mathbf{bPb}
 \end{array}$$

We construct the datatype *Pal* by interpreting the nonterminal as datatype and the names of the productions as names of the constructors:

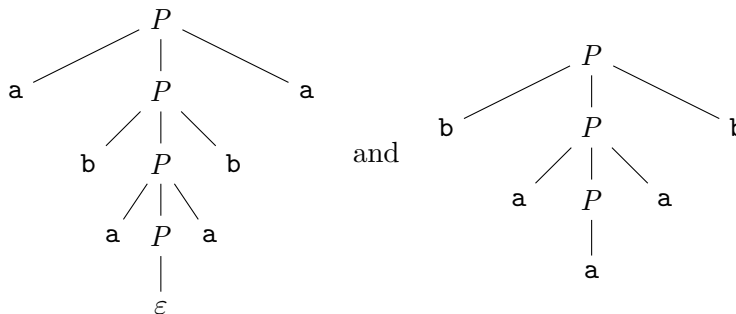
**data** *Pal* = *Empty* | *A* | *B* | *A*<sub>2</sub> *P* | *B*<sub>2</sub> *P*

Note that once again we keep the nonterminals on the right hand sides as arguments to the constructors, but omit all the terminal symbols.

The strings *abaaba* and *baaab* can be derived as follows:

$P \Rightarrow aPa \Rightarrow abPba \Rightarrow abaPaba \Rightarrow abaaba$   
 $P \Rightarrow bPb \Rightarrow baPba \Rightarrow baaab$

The parse trees corresponding to the derivations are the following:



Consequently, the desired Haskell definitions are

$pal_1 = A_2 (B_2 (A_2 \text{Empty}))$   
 $pal_2 = B_2 (A_2 A)$

## 2.22

1.

```
printPal :: Pal -> String
printPal Empty = ""
printPal A = "a"
printPal B = "b"
printPal (A2 p) = "a" ++ printPal p ++ "a"
printPal (B2 p) = "b" ++ printPal p ++ "b"
```

Note how the function follows the structure of the datatype *Pal* closely, and calls itself recursively wherever a recursive value of *Pal* occurs in the datatype. Such a pattern is typical for semantic functions.

B. Answers to exercises

2.

```
aCountPal :: Pal → Int
aCountPal Empty = 0
aCountPal A = 1
aCountPal B = 0
aCountPal (A2 p) = aCountPal p + 2
aCountPal (B2 p) = aCountPal p
```

**2.23** Recall the grammar from Exercise 2.8, this time with names for the productions:

```
MEmpty: M → ε
MA: M → aMa
MB: M → bMb
```

1. By systematically transforming the grammar, we obtain the following datatype:

```
data Mir = MEmpty | MA Mir | MB Mir
```

The concrete mirror palindromes  $cMir_1$  and  $cMir_2$  correspond to the following terms of type  $Mir$ :

```
aMir1 = MA (MB (MA Empty))
aMir2 = MA (MB (MB Empty))
```

2.

```
printMir :: Mir → String
printMir MEmpty = ""
printMir (MA m) = "a" ++ printMir m ++ "a"
printMir (MB m) = "b" ++ printMir m ++ "b"
```

3.

```
mirToPal :: Mir → Pal
mirToPal MEmpty = Empty
mirToPal (MA m) = A2 (mirToPal m)
mirToPal (MB m) = B2 (mirToPal m)
```

**2.24** Recall the grammar from Exercise 2.9, this time with names for the productions:

```
Stop: P → ε
POne: P → 1P1
PZeroL: P → 0P
PZeroR: P → P0
```

1.

```
data Parity = Stop | POne Parity | PZeroL Parity | PZeroR Parity
aEven1 = PZeroL (PZeroL (POne (PZeroL Stop)))
aEven2 = PZeroL (PZeroR (POne (PZeroL Stop)))
```

Note that the grammar is ambiguous, and other representations for  $cEven_1$  and  $cEven_2$  are possible, for instance:

```
aEven'1 = PZeroL (PZeroL (POne (PZeroR Stop)))
aEven'2 = PZeroR (PZeroL (POne (PZeroL Stop)))
```

2.

```
printParity :: Parity → String
printParity Stop = ""
printParity (POne p) = "1" ++ printParity p ++ "1"
printParity (PZeroL p) = "0" ++ printParity p
printParity (PZeroR p) = printParity p ++ "0"
```

**2.25** A grammar for bit lists that is not left-recursive is the following:

```
L → B Z | B
Z → , L Z | , L
B → 0 | 1
```

1.

```
data BitList = ConsBit Bit Z | SingleBit Bit
data Z = ConsBitList BitList Z | SingleBitList BitList
data Bit = Bit0 | Bit1
aBitList1 = ConsBit Bit0 (ConsBitList (SingleBit Bit1)
 (SingleBitList (SingleBit Bit0)))
aBitList2 = ConsBit Bit0 (ConsBitList (SingleBit Bit0)
 (SingleBitList (SingleBit Bit1)))
```

2.

```
printBitList :: BitList → String
printBitList (ConsBit b z) = printBit b ++ printZ z
printBitList (SingleBit b) = printBit b

printZ :: Z → String
printZ (ConsBitList bs z) = "," ++ printBitList bs ++ printZ z
printZ (SingleBitList bs) = "," ++ printBitList bs
```

B. Answers to exercises

```

printBit :: Bit → String
printBit Bit0 = "0"
printBit Bit1 = "1"

```

When multiple datatypes are involved, semantic functions typically still follow the structure of the datatypes closely. We get one function per datatype, and the functions call each other recursively where appropriate – we say they are *mutually recursive*.

3. We can still make the concatenation function structurally recursive in the first of the two bit lists. We never have to match on the second bit list:

```

concatBitList :: BitList → BitList → BitList
concatBitList (ConsBit b z) cs = ConsBit b (concatZ z cs)
concatBitList (SingleBit b) cs = ConsBit b (SingleBitList cs)
concatZ :: Z → BitList → Z
concatZ (ConsBitList bs z) cs = ConsBitList bs (concatZ z cs)
concatZ (SingleBitList bs) cs = ConsBitList bs (SingleBitList cs)

```

- 2.26** We only give the EBNF notation for the productions that change.

```

Digs → Dig*
Int → Sign? Nat
AlphaNums → AlphaNum*
AlphaNum → Letter | Dig

```

- 2.27**  $L(G?) = L(G) \cup \{\varepsilon\}$

**2.28**

1.  $L_1$  is generated by:

```

S → ZC
Z → aZb | ε
C → c*

```

and  $L_2$  is generated by:

```

S → AZ
A → a*
Z → bZc | ε

```

2. We have that

$$L_1 \cap L_2 = \{a^n b^n c^n \mid n \in \mathbb{N}\}$$

However, this language is not context-free, i. e., there is no context-free grammar that generates this language. We will see in Chapter 9 how to prove such a statement.

**2.29** No. Furthermore, for any language  $L$ , since  $\varepsilon \in L^*$ , we have  $\varepsilon \notin \overline{(L^*)}$ . On the other hand,  $\varepsilon \in (\overline{L})^*$ . Thus  $\overline{(L^*)}$  and  $(\overline{L})^*$  cannot be equal.

**2.30** For example  $L = \{x^n \mid n \in \mathbb{N}\}$ . For any language  $L$ , it holds that

$$L^* = (L^*)^*$$

so, given any language  $L$ , the language  $L^*$  fulfills the desired property.

**2.31** This is only the case when  $\varepsilon \notin L$ .

**2.32** No. the language  $L = \{\mathbf{aab}, \mathbf{baa}\}$  also satisfies  $L = L^R$ .

**2.33**

1. The shortest derivation is three steps long and yields the sentence **aa**. The sentences **baa**, **aba**, and **aab** can all be derived in four steps.
2. Several derivations are possible for the string **babbab**. Two of them are

$$\begin{aligned} \underline{S} &\Rightarrow \underline{AA} \Rightarrow \mathbf{bAA} \Rightarrow \mathbf{bAbA} \Rightarrow \mathbf{babA} \Rightarrow \mathbf{babbA} \Rightarrow \mathbf{babbAb} \Rightarrow \mathbf{babbab} \\ \underline{S} &\Rightarrow \underline{AA} \Rightarrow \underline{AAb} \Rightarrow \mathbf{bAAb} \Rightarrow \mathbf{bAbAb} \Rightarrow \mathbf{bAbbAb} \Rightarrow \mathbf{babbAb} \Rightarrow \mathbf{babbab} \end{aligned}$$

3. A leftmost derivation is:

$$\begin{aligned} \underline{S} &\Rightarrow \underline{AA} \Rightarrow^* \mathbf{b}^m \underline{AA} \Rightarrow^* \mathbf{b}^m \underline{Ab}^n A \Rightarrow \mathbf{b}^m \mathbf{ab}^n \underline{A} \Rightarrow^* \mathbf{b}^m \mathbf{ab}^n \underline{Ab}^p \\ &\Rightarrow \mathbf{b}^m \mathbf{ab}^n \mathbf{ab}^p \end{aligned}$$

**2.34** The grammar is equivalent to the grammar

$$\begin{aligned} S &\rightarrow \mathbf{aaB} \\ B &\rightarrow \mathbf{bBba} \\ B &\rightarrow \mathbf{a} \end{aligned}$$

This grammar generates the string **aaa** and the strings **aab<sup>m</sup>a(ba)<sup>m</sup>** for  $m \in \mathbb{N}$ ,  $m \geq 1$ . The string **aabbaabba** does not appear in this language.

**2.35** The language  $L$  is generated by:

$$S \rightarrow \mathbf{aS} \mathbf{a} \mid \mathbf{bS} \mathbf{b} \mid \mathbf{c}$$

The derivation is:

$$\underline{S} \Rightarrow \mathbf{a} \underline{S} \mathbf{a} \Rightarrow \mathbf{ab} \underline{S} \mathbf{ba} \Rightarrow \mathbf{abcba}$$

**2.36** The language generated by the grammar is

$$\{\mathbf{a}^n \mathbf{b}^n \mid n \in \mathbb{N}\}$$

B. Answers to exercises

The same language is also generated by the grammar

$$S \rightarrow \mathbf{aAb} \mid \varepsilon$$

**2.37** The first language is

$$\{\mathbf{a}^n \mid n \in \mathbb{N}\}$$

This language is also generated by the grammar

$$S \rightarrow \mathbf{aS} \mid \varepsilon$$

The second language is

$$\{\varepsilon\} \cup \{\mathbf{a}^{2n+1} \mid n \in \mathbb{N}\}$$

This language is also generated by the grammar

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow \mathbf{a} \mid \mathbf{aAa} \end{aligned}$$

or using EBNF notation

$$S \rightarrow (\mathbf{a(aa)^*})?$$

**2.38** All three grammars generate the language

$$\{\mathbf{a}^n \mid n \in \mathbb{N}\}$$

**2.39**

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow \mathbf{aAb} \mid \mathbf{ab} \end{aligned}$$

**2.40** The language is

$$L = \{\mathbf{a}^{2n+1} \mid n \in \mathbb{N}\}$$

A grammar for  $L$  without left-recursive productions is

$$A \rightarrow \mathbf{aaA} \mid \mathbf{a}$$

And a grammar without right-recursive productions is

$$A \rightarrow \mathbf{Aaa} \mid \mathbf{a}$$

**2.41** The language is



$$\{ab^n \mid n \in \mathbb{N}\}$$

A grammar for  $L$  without left-recursive productions is

$$\begin{aligned} X &\rightarrow aY \\ Y &\rightarrow bY \mid \varepsilon \end{aligned}$$

A grammar for  $L$  without left-recursive productions that is also non-contracting is

$$\begin{aligned} X &\rightarrow aY \mid a \\ Y &\rightarrow bY \mid b \end{aligned}$$

**2.42** A grammar that uses only productions with two or less symbols on the right hand side:

$$\begin{aligned} S &\rightarrow T \mid US \\ T &\rightarrow Xa \mid Ua \\ X &\rightarrow aS \\ U &\rightarrow S \mid YT \\ Y &\rightarrow SU \end{aligned}$$

The sentential forms  $aS$  and  $SU$  have been abstracted to nonterminals  $X$  and  $Y$ .

A grammar for the same language with only two nonterminals:

$$\begin{aligned} S &\rightarrow aSa \mid Ua \mid US \\ U &\rightarrow S \mid SUaSa \mid SUUa \end{aligned}$$

The nonterminal  $T$  has been substituted for its alternatives  $aSa$  and  $Ua$ .

**2.43**

$$\begin{aligned} S &\rightarrow 1O \\ O &\rightarrow 1O \mid 0N \\ N &\rightarrow 1^* \end{aligned}$$

**2.44** The language is generated by the grammar:

$$\begin{aligned} S &\rightarrow (A) \mid SS \\ A &\rightarrow S \mid \varepsilon \end{aligned}$$

A derivation for  $( ) ( ( ) ) ( )$  is:

$$\begin{aligned} \underline{S} &\Rightarrow \underline{SS} \Rightarrow \underline{SSS} \Rightarrow (\underline{A})SS \Rightarrow ()\underline{SS} \Rightarrow () (\underline{A})S \Rightarrow () (\underline{S})S \Rightarrow () ((\underline{A}))S \\ &\Rightarrow () (())\underline{S} \Rightarrow () (()) (\underline{A}) \Rightarrow () (()) () \end{aligned}$$

**2.45** The language is generated by the grammar

B. Answers to exercises

$$S \rightarrow (A) \mid [A] \mid SS$$

$$A \rightarrow S \mid \varepsilon$$

A derivation for  $[(\ )](\ )$  is:

$$\underline{S} \Rightarrow \underline{SS} \Rightarrow [\underline{A}]S \Rightarrow [\underline{S}]S \Rightarrow [(\underline{A})]S \Rightarrow [(\ )]\underline{S} \Rightarrow [(\ )](\underline{A}) \Rightarrow [(\ )](\ )$$

**2.46** First leftmost derivation:

$$\begin{aligned} & \underline{Sentence} \\ \Rightarrow & \underline{Subject} \underline{Predicate} \\ \Rightarrow & \mathbf{they} \underline{Predicate} \\ \Rightarrow & \mathbf{they} \underline{Verb} \underline{NounPhrase} \\ \Rightarrow & \mathbf{they} \mathbf{are} \underline{NounPhrase} \\ \Rightarrow & \mathbf{they} \mathbf{are} \underline{Adjective} \underline{Noun} \\ \Rightarrow & \mathbf{they} \mathbf{are} \mathbf{flying} \underline{Noun} \\ \Rightarrow & \mathbf{they} \mathbf{are} \mathbf{flying} \mathbf{planes} \end{aligned}$$

Second leftmost derivation:

$$\begin{aligned} & \underline{Sentence} \\ \Rightarrow & \underline{Subject} \underline{Predicate} \\ \Rightarrow & \mathbf{they} \underline{Predicate} \\ \Rightarrow & \mathbf{they} \underline{AuxVerb} \underline{Verb} \underline{Noun} \\ \Rightarrow & \mathbf{they} \mathbf{are} \underline{Verb} \underline{Noun} \\ \Rightarrow & \mathbf{they} \mathbf{are} \mathbf{flying} \underline{Noun} \\ \Rightarrow & \mathbf{they} \mathbf{are} \mathbf{flying} \mathbf{planes} \end{aligned}$$

**2.48** Here is an unambiguous grammar for the language from Exercise 2.45:

$$S \rightarrow (E)E \mid [E]E$$

$$E \rightarrow \varepsilon \mid S$$

**2.49** Here is a leftmost derivation for  $\clubsuit \diamond \clubsuit \triangle \spadesuit$ .

$$\begin{aligned} \odot & \Rightarrow \odot \triangle \otimes \Rightarrow \otimes \triangle \otimes \Rightarrow \otimes \diamond \oplus \triangle \otimes \Rightarrow \oplus \diamond \oplus \triangle \otimes \Rightarrow \clubsuit \diamond \oplus \triangle \otimes \\ & \Rightarrow \clubsuit \diamond \clubsuit \triangle \otimes \Rightarrow \clubsuit \diamond \clubsuit \triangle \oplus \Rightarrow \clubsuit \diamond \clubsuit \triangle \spadesuit \end{aligned}$$

Notice that the grammar of this exercise is the same, up to renaming, as the grammar

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow 0 \mid 1$$

**2.50** The palindrome  $\varepsilon$  with length 0 can be generated by the grammar with the derivation  $P \Rightarrow \varepsilon$ . The palindromes **a**, **b**, and **c** are the palindromes with length 1. They can be derived with  $P \Rightarrow \mathbf{a}$ ,  $P \Rightarrow \mathbf{b}$ ,  $P \Rightarrow \mathbf{c}$ , respectively.

Suppose that the palindrome  $s$  with a length of 2 or more. Then  $s$  can be written as **ata** or **btb** or **ctc** where the length of  $t$  is strictly smaller, and  $t$  also is a palindrome. Thus, by induction hypothesis, there is a derivation  $P \Rightarrow^* t$ . But then, there is also a derivation for  $s$ , for example  $P \Rightarrow^* t \Rightarrow \mathbf{ata}$  in the first situation – the other two cases are analogous.

We have now proved that any palindrome can be generated by the grammar – we still have to prove that anything generated by the grammar is a palindrome, but this is easy to see by induction over the length of derivations. Certainly  $\varepsilon$ ,  $a$ ,  $b$ , and  $c$  are palindromes. And if  $s$  is a palindrome that can be derived, so are **asa**, **bsb**, and **csc**.

**3.1** Either we use the predefined predicate *isUpper* in module *Data.Char*,

$$\mathit{capital} = \mathit{satisfy} \ \mathit{isUpper}$$

or we make use of the ordering defined characters,

$$\mathit{capital} = \mathit{satisfy} \ (\lambda s \rightarrow ('A' \leq s) \wedge (s \leq 'Z'))$$

**3.2** A symbol equal to  $a$  satisfies the predicate  $(= a)$ :

$$\mathit{symbol} \ a = \mathit{satisfy} \ (= a)$$

**3.3** The function *epsilon* is a special case of *succeed*:

$$\begin{aligned} \mathit{epsilon} &:: \mathit{Parser} \ s \ () \\ \mathit{epsilon} &= \mathit{succeed} \ () \end{aligned}$$

**3.4** Let  $xs :: [s]$ . Then

$$\begin{aligned} &(f \langle \$ \rangle \mathit{succeed} \ a) \ xs \\ = &\{ \text{definition of } \langle \$ \rangle \} \\ &[(f \ x, \ ys) \mid (x, \ ys) \leftarrow \mathit{succeed} \ a \ xs] \\ = &\{ \text{definition of } \mathit{succeed} \} \\ &[(f \ a, \ xs)] \\ = &\{ \text{definition of } \mathit{succeed} \} \\ &\mathit{succeed} \ (f \ a) \ xs \end{aligned}$$

**3.5** The type and results of  $(: \langle \$ \rangle \mathit{symbol} \ 'a')$  are (note that you cannot write this as a definition in Haskell):

B. Answers to exercises

```
((: <$> symbol 'a') :: Parser Char (String → String)
(: <$> symbol 'a') [] = []
(: <$> symbol 'a') (x : xs) | x == 'a' = [(x :), xs]
| otherwise = []
```

**3.6** The type and results of  $(: <$> \text{symbol 'a'} <*> p)$  are:

```
((: <$> symbol 'a' <*> p) :: Parser Char String
(: <$> symbol 'a' <*> p) [] = []
(: <$> symbol 'a' <*> p) (x : xs)
| x == 'a' = [('a' : x, ys) | (x, ys) ← p xs]
| otherwise = []
```

**3.7**

```
pBool :: Parser Char Bool
pBool = const True <$> token "True"
<|> const False <$> token "False"
```

**3.9**

1.

```
data Pal2 = Nil | Leafa | Leafb | Twa Pal2 | Twob Pal2
```

2.

```
palin2 :: Parser Char Pal2
palin2 = (_ y _ → Twa y) <$>
 symbol 'a' <*> palin2 <*> symbol 'a'
<|> (_ y _ → Twob y) <$>
 symbol 'b' <*> palin2 <*> symbol 'b'
<|> const Leafa <$> symbol 'a'
<|> const Leafb <$> symbol 'b'
<|> succeed Nil
```

3.

```
palina :: Parser Char Int
palina = (_ y _ → y + 2) <$>
 symbol 'a' <*> palina <*> symbol 'a'
<|> (_ y _ → y) <$>
 symbol 'b' <*> palina <*> symbol 'b'
<|> const 1 <$> symbol 'a'
<|> const 0 <$> symbol 'b'
<|> succeed 0
```

### 3.10

1.

```

data English = E1 Subject Pred
data Subject = E2 String
data Pred = E3 Verb NounP | E4 AuxV Verb Noun
data Verb = E5 String | E6 String
data AuxV = E7 String
data NounP = E8 Adj Noun
data Adj = E9 String
data Noun = E10 String

```

2.

```

english :: Parser Char English
english = E1 <$> subject <*> pred
subject = E2 <$> token "they"
pred = E3 <$> verb <*> nounp
 <|> E4 <$> auxv <*> verb <*> noun
verb = E5 <$> token "are"
 <|> E6 <$> token "flying"
auxv = E7 <$> token "are"
nounp = E8 <$> adj <*> noun
adj = E9 <$> token "flying"
noun = E10 <$> token "planes"

```

**3.11** As <|> uses  $\#$ , it is more efficiently evaluated if right-associative.

**3.12** The function is the same as <\*>, but instead of applying the result of the first parser to that of the second, it pairs them together:

$$\begin{aligned}
(<, >) &:: \text{Parser } s \ a \rightarrow \text{Parser } s \ b \rightarrow \text{Parser } s \ (a, b) \\
(p <, > q) \ xs &= [((x, y), zs) \\
&\quad | (x, ys) \leftarrow p \ xs \\
&\quad , (y, zs) \leftarrow q \ ys \\
&\quad ]
\end{aligned}$$

**3.13** ‘Parser transformer’, or ‘parser modifier’ or ‘parser postprocessor’, etcetera.

**3.14** The transformer <\$> does to the result part of parsers what *map* does to the elements of a list.

**3.15** The parser combinators <\*> and <, > can be defined in terms of each other:

$$\begin{aligned}
p <*> q &= \text{uncurry } (\$) <$> (p <, > q) \\
p <, > q &= (,) <$> p <*> q
\end{aligned}$$



The definition of the parser *sumParser* is:

```
sumParser :: Parser Char Int
sumParser = chainr newdigit plusParser
```

```
?> sumParser "1+2+3"
[(6, ""), (3, "+3"), (1, "+2+3")]
?> sumParser "1+2+a"
[(3, "+a"), (1, "+2+a")]
?> sumParser "1"
[(1, "")]
```

Note that the parser also recognises a single integer.

3. The parser *many* should be replaced by the parser *greedy* in the definition of *listOf*.

**3.21** We introduce the abbreviation

$$listOfa = (:) \langle \$ \rangle \textit{symbol} \textit{'a'}$$

and use the results of Exercises 3.5 and 3.6.

$xs = []$ :

$$\begin{aligned} & \textit{many} (\textit{symbol} \textit{'a'}) [] \\ = & \{ \textit{definition of } \textit{many} \textit{ and } \textit{listOfa} \} \\ & (\textit{listOfa} \langle * \rangle \textit{many} (\textit{symbol} \textit{'a'}) \langle | \rangle \textit{succeed} []) [] \\ = & \{ \textit{definition of } \langle | \rangle \} \\ & (\textit{listOfa} \langle * \rangle \textit{many} (\textit{symbol} \textit{'a'})) [] \textit{++} \textit{succeed} [] [] \\ = & \{ \text{Exercise 3.6, definition of } \textit{succeed} \} \\ & [] \textit{++} [( [], [] )] \\ = & \{ \textit{definition of } \textit{++} \} \\ & [( [], [] )] \end{aligned}$$

$xs = [\textit{'a'}]$ :

$$\begin{aligned} & \textit{many} (\textit{symbol} \textit{'a'}) [\textit{'a'}] \\ = & \{ \textit{definition of } \textit{many} \textit{ and } \textit{listOfa} \} \\ & (\textit{listOfa} \langle * \rangle \textit{many} (\textit{symbol} \textit{'a'}) \langle | \rangle \textit{succeed} []) [\textit{'a'}] \\ = & \{ \textit{definition of } \langle | \rangle \} \\ & (\textit{listOfa} \langle * \rangle \textit{many} (\textit{symbol} \textit{'a'})) [\textit{'a'}] \textit{++} \textit{succeed} [] [\textit{'a'}] \\ = & \{ \text{Exercise 3.6, previous calculation} \} \end{aligned}$$

B. Answers to exercises

$$[[('a'), []], ([], ['a'])]$$

$xs = ['b']$ :

$$\begin{aligned} & \text{many (symbol 'a')} ['b'] \\ = & \{ \text{as before} \} \\ & (\text{listOfa} \langle * \rangle \text{many (symbol 'a')}) ['b'] \# \text{succeed } [] ['b'] \\ = & \{ \text{Exercise 3.6, previous calculation} \} \\ & ([], ['b']) \end{aligned}$$

$xs = ['a', 'b']$ :

$$\begin{aligned} & \text{many (symbol 'a')} ['a', 'b'] \\ = & \{ \text{as before} \} \\ & (\text{listOfa} \langle * \rangle \text{many (symbol 'a')}) ['a', 'b'] \# \text{succeed } [] ['a', 'b'] \\ = & \{ \text{Exercise 3.6, previous calculation} \} \\ & ([('a'), ['b']], ([], ['a', 'b'])) \end{aligned}$$

$xs = ['a', 'a', 'b']$ :

$$\begin{aligned} & \text{many (symbol 'a')} ['a', 'a', 'b'] \\ = & \{ \text{as before} \} \\ & (\text{listOfa} \langle * \rangle \text{many (symbol 'a')}) ['a', 'a', 'b'] \# \text{succeed } [] ['a', 'a', 'b'] \\ = & \{ \text{Exercise 3.6, previous calculation} \} \\ & ([('a', 'a'), ['b']], ([('a'), ['a', 'b']], ([], ['a', 'a', 'b']))) \end{aligned}$$

**3.22** The empty alternative is presented last, because the  $\langle | \rangle$  combinator uses list concatenation for concatenating lists of successes. This also holds for the recursive calls; thus the ‘greedy’ parsing of all three a’s is presented first, then two a’s with a singleton rest string, then one a, and finally the empty result with the original input as rest string.

**3.24**

— Combinators for repetition

$$\begin{aligned} \text{psequence} & :: [\text{Parser } s \ a] \rightarrow \text{Parser } s \ [a] \\ \text{psequence } [] & = \text{succeed } [] \\ \text{psequence } (p : ps) & = (:) \langle \$ \rangle p \langle * \rangle \text{psequence } ps \\ \text{psequence}' & :: [\text{Parser } s \ a] \rightarrow \text{Parser } s \ [a] \\ \text{psequence}' & = \text{foldr } f \ (\text{succeed } []) \\ & \text{where } f \ p \ q = (:) \langle \$ \rangle p \langle * \rangle q \end{aligned}$$



```
choice :: [Parser s a] → Parser s a
choice = foldr (<|>) failp
```

```
?> (psequence [digit, satisfy isUpper]) "1A"
[("1A","")]
?> (psequence [digit, satisfy isUpper]) "1Ab"
[("1A","b")]
?> (psequence [digit, satisfy isUpper]) "1ab"
[]
```

```
?> (choice [digit, satisfy isUpper]) "1ab"
[('1',"ab")]
?> (choice [digit, satisfy isUpper]) "Ab"
[('A',"b")]
?> (choice [digit, satisfy isUpper]) "ab"
[]
```

### 3.25

```
token :: Eq s ⇒ [s] → Parser s [s]
token = psequence . map symbol
```

### 3.27

```
identifier :: Parser Char String
identifier = (:) <$> satisfy isAlpha <*> greedy (satisfy isAlphaNum)
```

### 3.28

1. As Haskell terms:

```
"abc": Var "abc"
"(abc)": Var "abc"
"a*b+1": Var "a" :* Var "b" :+: Con 1
"a*(b+1)": Var "a" :* (Var "b" :+: Con 1)
"-1-a": Con (-1) :-: Var "a"
"a(1,b)": Fun "a" [Con 1, Var "b"]
```

2. The parser *fact* first tries to parse an integer, then a variable, then a function application and finally a parenthesised expression. A function application is a variable followed by an argument list. When the parser encounters a function application, a variable will first be recognised. This first solution will however

## B. Answers to exercises

not lead to a parse tree for the complete expression because the list of arguments that comes after the variable cannot be parsed.

If we swap the second and the third line in the definition of the parser *fact*, the parse tree for a function application will be the first solution of the parser:

```
fact :: Parser Char Expr
fact = Con <$> integer
 <|> Fun <$> identifier <*> parenthesised (commaList expr)
 <|> Var <$> identifier
 <|> parenthesised expr
```

```
?> expr "a(1,b)"
[(Fun "a" [Con 1,Var "b"],""), (Var "a","(1,b)")]
```

**3.29** A function with no arguments is not accepted by the parser:

```
?> expr "f()"
[(Var "f","()")]
```

The parser *parenthesised (commaList expr)* that is used in the parser *fact* does not accept an empty list of arguments because *commaList* does not. To accept an empty list we modify the parser *fact* as follows:

```
fact :: Parser Char Expr
fact = Con <$> integer
 <|> Fun <$> identifier
 <*> parenthesised (commaList expr <|> succeed [])
 <|> Var <$> identifier
 <|> parenthesised expr
```

```
?> expr "f()"
[(Fun "f" [],""), (Var "f","()")]
```

**3.30**

```
expr = chainr (chainl term (const (:-) <$> symbol '-'))
 (const (:+) <$> symbol '+')
```

**3.31** The datatype *Expr* is extended as follows to allow raising an expression to the power of an expression:

```
data Expr = Con Int
 | Var String
```

```

| Fun String [Expr]
| Expr :+: Expr
| Expr :-: Expr
| Expr **: Expr
| Expr :/: Expr
| Expr : ^ : Expr

```

**deriving** Show

Now the parser  $expr'$  of Listing 3.8 can be extended with a new level of priorities:

```

powis = [('^ ', (: ^ :))]
expr' :: Parser Char Expr
expr' = foldr gen fact' [addis, multis, powis]

```

Note that because of the use of *chainl* all the operators listed in *addis*, *multis* and *powis* are treated as left-associative.

**3.32** The proofs can be given by using laws for list comprehension, but here we prefer to exploit the following equation

$$(f \langle \$ \rangle p) \, xs = \text{map } (f \, ** \, id) \, (p \, xs) \quad (\text{B.1})$$

where  $(**)$  is defined by

$$\begin{aligned} (** ) &:: (a \rightarrow c) \rightarrow (b \rightarrow d) \rightarrow (a, b) \rightarrow (c, d) \\ (f \, ** \, g) \, (a, b) &= (f \, a, g \, b) \end{aligned}$$

It has the following property:

$$(f \, ** \, g) \cdot (h \, ** \, k) = (f \cdot h) \, ** \, (g \cdot k) \quad (\text{B.2})$$

Furthermore, we will use the following laws about *map* in our proof: *map* distributes over composition, concatenation, and the function *concat*:

$$\text{map } f \cdot \text{map } g = \text{map } (f \cdot g) \quad (\text{B.3})$$

$$\text{map } f \, (x \, \# \, y) = \text{map } f \, x \, \# \, \text{map } f \, y \quad (\text{B.4})$$

$$\text{map } f \cdot \text{concat} = \text{concat} \cdot \text{map } (\text{map } f) \quad (\text{B.5})$$

1.

$$\begin{aligned} & (h \langle \$ \rangle (f \langle \$ \rangle p)) \, xs \\ &= \{ (\text{B.1}) \} \\ & \quad \text{map } (h \, ** \, id) \, ((f \langle \$ \rangle p) \, xs) \\ &= \{ (\text{B.1}) \} \\ & \quad \text{map } (h \, ** \, id) \, (\text{map } (f \, ** \, id) \, (p \, xs)) \end{aligned}$$

B. Answers to exercises

$$\begin{aligned}
&= \{ \text{(B.3)} \} \\
&\quad \text{map } ((h \text{ ** } id) . (f \text{ ** } id)) (p \text{ } xs) \\
&= \{ \text{(B.2)} \} \\
&\quad \text{map } ((h . f) \text{ ** } id) (p \text{ } xs) \\
&= \{ \text{(B.1)} \} \\
&\quad ((h . f) \langle \$ \rangle p) \text{ } xs
\end{aligned}$$

2.

$$\begin{aligned}
&(h \langle \$ \rangle (p \langle | \rangle q)) \text{ } xs \\
&= \{ \text{(B.1)} \} \\
&\quad \text{map } (h \text{ ** } id) ((p \langle | \rangle q) \text{ } xs) \\
&= \{ \text{definition of } \langle | \rangle \} \\
&\quad \text{map } (h \text{ ** } id) (p \text{ } xs \text{ } ++ \text{ } q \text{ } xs) \\
&= \{ \text{(B.4)} \} \\
&\quad \text{map } (h \text{ ** } id) (p \text{ } xs) \text{ } ++ \text{ } \text{map } (h \text{ ** } id) (q \text{ } xs) \\
&= \{ \text{(B.1)} \} \\
&\quad (h \langle \$ \rangle p) \text{ } xs \text{ } ++ (h \langle \$ \rangle q) \text{ } xs \\
&= \{ \text{definition of } \langle | \rangle \} \\
&\quad ((h \langle \$ \rangle p) \langle | \rangle (h \langle \$ \rangle q)) \text{ } xs
\end{aligned}$$

3. First note that  $(p \langle * \rangle q) \text{ } xs$  can be written as

$$(p \langle * \rangle q) \text{ } xs = \text{concat } (\text{map } (mc \text{ } q) (p \text{ } xs)) \tag{B.6}$$

where

$$mc \text{ } q (f, ys) = \text{map } (f \text{ ** } id) (q \text{ } ys)$$

Now we calculate

$$\begin{aligned}
&(((h.) \langle \$ \rangle p) \langle * \rangle q) \text{ } xs \\
&= \{ \text{(B.6)} \} \\
&\quad \text{concat } (\text{map } (mc \text{ } q) (((h.) \langle \$ \rangle p) \text{ } xs)) \\
&= \{ \text{(B.1)} \} \\
&\quad \text{concat } (\text{map } (mc \text{ } q) (\text{map } ((h.) \text{ ** } id) (p \text{ } xs))) \\
&= \{ \text{(B.3)} \} \\
&\quad \text{concat } (\text{map } (\text{map } (h \text{ ** } id) (\text{map } (mc \text{ } q) (p \text{ } xs)))) \\
&= \{ \text{(B.7), see below} \} \\
&\quad \text{concat } (\text{map } ((\text{map } (h \text{ ** } id)) . mc \text{ } q) (p \text{ } xs))
\end{aligned}$$

$$\begin{aligned}
&= \{ \text{(B.3)} \} \\
&\quad \text{concat} (\text{map} (\text{map} (h \text{**} id)) (\text{map} (mc\ q) (p\ xs))) \\
&= \{ \text{(B.5)} \} \\
&\quad \text{map} (h \text{**} id) (\text{concat} (\text{map} (mc\ q) (p\ xs))) \\
&= \{ \text{(B.6)} \} \\
&\quad \text{map} (h \text{**} id) ((p \langle * \rangle q)\ xs) \\
&= \{ \text{(B.1)} \} \\
&\quad (h \langle \$ \rangle (p \langle * \rangle q))\ xs
\end{aligned}$$

It remains to prove the claim

$$mc\ q . ((h.) \text{**} id) = \text{map} (h \text{**} id) . mc\ q \tag{B.7}$$

This claim is also proved by calculation:

$$\begin{aligned}
&((\text{map} (h \text{**} id)) . mc\ q) (f, ys) \\
&= \{ \text{definition of } . \} \\
&\quad \text{map} (h \text{**} id) (mc\ q (f, ys)) \\
&= \{ \text{definition of } mc\ q \} \\
&\quad \text{map} (h \text{**} id) (\text{map} (f \text{**} id) (q\ ys)) \\
&= \{ \text{map and ** distribute over composition} \} \\
&\quad \text{map} ((h . f) \text{**} id) (q\ ys) \\
&= \{ \text{definition of } mc\ q \} \\
&\quad mc\ q (h . f, ys) \\
&= \{ \text{definition of **} \} \\
&\quad (mc\ q . ((h.) \text{**} id)) (f, ys)
\end{aligned}$$

### 3.33

```

pMir :: Parser Char Mir
pMir = (_ m _ → MB m) <$> symbol 'b' <*> pMir <*> symbol 'b'
 <|> (_ m _ → MA m) <$> symbol 'a' <*> pMir <*> symbol 'a'
 <|> succeed MEmpty

```

### 3.34

```

pBitList :: Parser Char BitList
pBitList = SingleB <$> pBit
 <|> (\b _ bs → ConsB b bs) <$> pBit <*> symbol ',' <*> pBitList
pBit = const Bit0 <$> symbol '0'
 <|> const Bit1 <$> symbol '1'

```

**3.35**

— Parser for floating point numbers

```
fixed :: Parser Char Float
fixed = (+) <$> (fromIntegral <$> greedy integer)
 <*> (((λ_ y → y) <$> symbol ' . ' <*> fractpart) 'option' 0.0)

fractpart :: Parser Char Float
fractpart = foldr f 0.0 <$> greedy newdigit
 where f d n = (n + fromIntegral d) / 10.0
```

**3.36**

```
float :: Parser Char Float
float = f <$> fixed
 <*> (((λ_ y → y) <$> symbol 'E' <*> integer) 'option' 0)
 where f m e = m * power e
 power e | e < 0 = 1.0 / power (-e)
 | otherwise = fromIntegral (10e)
```

**3.37** Parse trees for Java assignments are of type:

```
data JavaAssign = JAssign String Expr
 deriving Show
```

The parser is defined as follows:

```
assign :: Parser Char JavaAssign
assign = JAssign
 <$> identifier
 <*> ((λ_ y → y) <$> symbol '=' <*> expr <*> symbol ';' ;)
```

```
?> assign "x1=(a+1)*2;"
[(JAssign "x1" (Var "a" :+: Con 1 :* Con 2), "")]
?> assign "x=a+1"
[]
```

Note that the second example is not recognised as an assignment because the string does not end with a semicolon.

**4.1**

```
data FloatLiteral = FL1 IntPart FractPart ExponentPart FloatSuffix
 | FL2 FractPart ExponentPart FloatSuffix
 | FL3 IntPart ExponentPart FloatSuffix
```

```

 | FL4 IntPart ExponentPart FloatSuffix
deriving Show
type ExponentPart = String
type ExponentIndicator = String
type SignedInteger = String
type IntPart = String
type FractPart = String
type FloatSuffix = String

digit = satisfy isDigit
digits = many1 digit
floatLiteral =
 (λ a b c d e → FL1 a c d e)
 <$> intPart <*> period <*> optfract <*> optexp <*> optfloat
<|> (λ a b c d → FL2 b c d)
 <$> period <*> fractPart <*> optexp <*> optfloat
<|> (λ a b c → FL3 a b c)
 <$> intPart <*> exponentPart <*> optfloat
<|> (λ a b c → FL4 a b c)
 <$> intPart <*> optexp <*> floatSuffix

intPart = signedInteger
fractPart = digits
exponentPart = (+) <$> exponentIndicator <*> signedInteger
signedInteger = (+) <$> option sign "" <*> digits
exponentIndicator = token "e" <|> token "E"
sign = token "+" <|> token "-"
floatSuffix = token "f" <|> token "F"
 <|> token "d" <|> token "D"
period = token "."
optexp = option exponentPart ""
optfract = option fractPart ""
optfloat = option floatSuffix ""

```

**4.2** The data and type definitions are the same as before, only the parsers return another (semantic) result.

```

digit = f <$> satisfy isDigit
 where f c = ord c - ord '0'
digits = foldl f 0 <$> many1 digit
 where f a b = 10*a + b
floatLiteral =
 (λ a b c d e → (fromIntegral a + c) * power d) <$>
 intPart <*> period <*> optfract <*> optexp <*> optfloat
<|> (λ a b c d → b * power c) <$>
 period <*> fractPart <*> optexp <*> optfloat

```

B. Answers to exercises

```

<|> (\a b c -> (fromIntegral a) * power b) <$>
 intPart <*> exponentPart <*> optfloat
<|> (\a b c -> (fromIntegral a) * power b) <$>
 intPart <*> optexp <*> floatSuffix
intPart = signedInteger
fractPart = foldr f 0.0 <$> many1 digit
 where f a b = (fromIntegral a + b)/10
exponentPart = (\x y -> y) <$> exponentIndicator <*> signedInteger
signedInteger = (\ x y -> x y) <$> option sign id <*> digits
exponentIndicator = symbol 'e' <|> symbol 'E'
sign = const id <$> symbol '+'
 <|> const negate <$> symbol '-'
floatSuffix = symbol 'f' <|> symbol 'F' <|> symbol 'd' <|> symbol 'D'
period = symbol '.'
optexp = option exponentPart 0
optfract = option fractPart 0.0
optfloat = option floatSuffix ''
power e | e < 0 = 1 / power (-e)
 | otherwise = fromIntegral (10^e)

```

4.3 The parsing scheme for Java floats is

```

digit = f <$> satisfy isDigit
 where f c =
digits = f <$> many1 digit
 where f ds = ..
floatLiteral = f1 <$>
 intPart <*> period <*> optfract <*> optexp <*> optfloat
<|> f2 <$>
 period <*> fractPart <*> optexp <*> optfloat
<|> f3 <$>
 intPart <*> exponentPart <*> optfloat
<|> f4 <$>
 intPart <*> optexp <*> floatSuffix
 where
 f1 a b c d e =
 f2 a b c d =
 f3 a b c =
 f4 a b c =
intPart = signedInteger
fractPart = f <$> many1 digit
 where f ds =
exponentPart = f <$> exponentIndicator <*> signedInteger
 where f x y =

```



```

signedInteger = f <$> option sign ?? <*> digits
 where f x y =
exponentIndicator = f1 <$> symbol 'e' <|> f2 <$> symbol 'E'
 where
 f1 c =
 f2 c = ..
sign = f1 <$> symbol '+' <|> f2 <$> symbol '-'
 where
 f1 h =
 f2 h =
floatSuffix = f1 <$> symbol 'f'
 <|> f2 <$> symbol 'F'
 <|> f3 <$> symbol 'd'
 <|> f4 <$> symbol 'D'
 where
 f1 c =
 f2 c =
 f3 c =
 f4 c =
period = symbol '.'
optexp = option exponentPart ??
optfract = option fractPart ??
optfloat = option floatSuffix ??

```

## 5.1

```

type LNTreeAlgebra a b x = (a → x, x → b → x → x)
foldLNTree :: LNTreeAlgebra a b x → LNTree a b → x
foldLNTree (leaf, node) = fold
 where
 fold (Leaf a) = leaf a
 fold (Node l m r) = node (fold l) m (fold r)

```

## 5.2

1. Definition of *height* by case analysis:

$$\begin{aligned}
 \text{height } (\text{Leaf } x) &= 0 \\
 \text{height } (\text{Bin } lt \text{ } rt) &= 1 + (\text{height } lt \text{ 'max' } \text{height } rt)
 \end{aligned}$$

Definition as a fold:

$$\begin{aligned}
 \text{height} &:: \text{BinTree } x \rightarrow \text{Int} \\
 \text{height} &= \text{foldBinTree heightAlgebra} \\
 \text{heightAlgebra} &= (\lambda u \ v \rightarrow 1 + (u \text{ 'max' } v), \text{const } 0)
 \end{aligned}$$

## B. Answers to exercises

2. Definition of *flatten* by case analysis:

$$\begin{aligned} \text{flatten } (\text{Leaf } x) &= [x] \\ \text{flatten } (\text{Bin } lt \ rt) &= \text{flatten } lt \ ++ \ \text{flatten } rt \end{aligned}$$

Definition as a fold:

$$\begin{aligned} \text{flatten} &:: \text{BinTree } x \rightarrow [x] \\ \text{flatten} &= \text{foldBinTree } ((\ ++ \ ), \lambda x \rightarrow [x]) \end{aligned}$$

3. Definition of *maxBinTree* by case analysis:

$$\begin{aligned} \text{maxBinTree } (\text{Leaf } x) &= x \\ \text{maxBinTree } (\text{Bin } lt \ rt) &= \text{maxBinTree } lt \ \text{'max'} \ \text{maxBinTree } rt \end{aligned}$$

Definition as a fold:

$$\begin{aligned} \text{maxBinTree} &:: \text{Ord } x \Rightarrow \text{BinTree } x \rightarrow x \\ \text{maxBinTree} &= \text{foldBinTree } (\text{max}, \text{id}) \end{aligned}$$

4. Definition of *sp* by case analysis:

$$\begin{aligned} \text{sp } (\text{Leaf } x) &= 0 \\ \text{sp } (\text{Bin } lt \ rt) &= 1 + (\text{sp } lt) \ \text{'min'} \ (\text{sp } rt) \end{aligned}$$

Definition as a fold:

$$\begin{aligned} \text{sp} &:: \text{BinTree } x \rightarrow \text{Int} \\ \text{sp} &= \text{foldBinTree } \text{spAlgebra} \\ \text{spAlgebra} &= (\lambda u \ v \rightarrow 1 + u \ \text{'min'} \ v, \text{const } 0) \end{aligned}$$

5. Definition of *mapBinTree* by case analysis:

$$\begin{aligned} \text{mapBinTree } f \ (\text{Leaf } x) &= \text{Leaf } (f \ x) \\ \text{mapBinTree } f \ (\text{Bin } lt \ rt) &= \text{Bin } (\text{mapBinTree } f \ lt) \ (\text{mapBinTree } f \ rt) \end{aligned}$$

Definition as a fold:

$$\begin{aligned} \text{mapBinTree} &:: (a \rightarrow b) \rightarrow \text{BinTree } a \rightarrow \text{BinTree } b \\ \text{mapBinTree } f &= \text{foldBinTree } (\text{Bin}, \text{Leaf } . f) \end{aligned}$$

### 5.3 Using explicit recursion:

$$\begin{aligned} \text{allPaths } (\text{Leaf } x) &= [[]] \\ \text{allPaths } (\text{Bin } lt \ rt) &= \text{map } (L:) \ (\text{allPaths } lt) \\ &\quad \ ++ \ \text{map } (R:) \ (\text{allPaths } rt) \end{aligned}$$

As a fold:

```
allPaths :: BinTree a → [[Direction]]
allPaths = foldBinTree psAlgebra
psAlgebra :: BinTreeAlgebra a [[Direction]]
psAlgebra = (λu v → map (L:) u ++ map (R:) v, const [[]])
```

## 5.4

1.

```
data Resist = Resist |: Resist
 | Resist *: Resist
 | BasicR Float
 deriving Show
type ResistAlgebra a = (a → a → a, a → a → a, Float → a)
foldResist :: ResistAlgebra a → Resist → a
foldResist (par, seq, basic) = fold
 where
 fold (r1 |: r2) = par (fold r1) (fold r2)
 fold (r1 *: r2) = seq (fold r1) (fold r2)
 fold (BasicR f) = basic f
```

2.

```
result :: Resist → Float
result = foldResist resultAlgebra
resultAlgebra :: ResistAlgebra Float
resultAlgebra = (λu v → (u * v) / (u + v), (+), id)
```

## 5.5

1. isSum = foldExpr ((&&)  
                    , \x y -> False  
                    , \x y -> False  
                    , \x y -> False  
                    , const True  
                    , const True  
                    , \x -> (&&)  
                    )
2. vars = foldExpr ((++)  
                    , (++)  
                    , (++)

B. Answers to exercises

```

, (++)
, const []
, \x -> [x]
, \x y z -> x : (y ++ z)
)

```

5.6

1. *der* computes a symbolic differentiation of an expression.
2. The function *der* is not a compositional function on *Expr* because the righthand sides of the 'Mul' and 'Dvd' expressions do not only use *der e<sub>1</sub> dx* and *der e<sub>2</sub> dx*, but also *e<sub>1</sub>* and *e<sub>2</sub>* themselves.
- 3.

```

data Exp = Exp 'Plus' Exp
 | Exp 'Sub' Exp
 | Con Float
 | Idf String
deriving Show
type ExpAlgebra a = (a -> a -> a
 , a -> a -> a
 , Float -> a
 , String -> a
)
foldExp :: ExpAlgebra a -> Exp -> a
foldExp (plus, sub, con, idf) = fold
where
 fold (e1 'Plus' e2) = plus (fold e1) (fold e2)
 fold (e1 'Sub' e2) = sub (fold e1) (fold e2)
 fold (Con n) = con n
 fold (Idf s) = idf s

```

4. Using explicit recursion:

```

der :: Exp -> String -> Exp
der (e1 'Plus' e2) dx = der e1 dx 'Plus' der e2 dx
der (e1 'Sub' e2) dx = der e1 dx 'Sub' der e2 dx
der (Con f) dx = Con 0
der (Idf s) dx = if s == dx then Con 1 else Con 0

```

Using a fold:

```

der = foldExp derAlgebra
derAlgebra :: ExpAlgebra (String -> Exp)

```

```

derAlgebra = (λf g → λs → f s 'Plus' g s
 , λf g → λs → f s 'Sub' g s
 , λn → λs → Con 0
 , λs → λt → if s == t then Con 1 else Con 0
)

```

5.7 Using explicit recursion:

```

replace (Leaf x) y = Leaf y
replace (Bin lt rt) y = Bin (replace lt y) (replace rt y)

```

As a fold:

```

replace :: BinTree a → a → BinTree a
replace = foldBinTree repAlgebra
repAlgebra = (λf g → λy → Bin (f y) (g y), λx y → Leaf y)

```

5.8 Using explicit recursion:

```

path2Value (Leaf x) =
 λbs → if null bs then x else error "no rootpath"
path2Value (Bin lt rt) =
 λbs → case bs of
 [] → error "no rootpath"
 (L : rs) → path2Value lt rs
 (R : rs) → path2Value rt rs

```

Using a fold:

```

path2Value :: BinTree a → [Direction] → a
path2Value = foldBinTree pvAlgebra
pvAlgebra :: BinTreeAlgebra a ([Direction] → a)
pvAlgebra = (λfl fr → λbs → case bs of
 [] → error "no rootpath"
 (L : rs) → fl rs
 (R : rs) → fr rs
 , λx → λbs → if null bs then x
 else error "no rootpath")

```

5.9

1.

```

type PalAlgebra p = (p, p, p, p → p, p → p)

```

B. Answers to exercises

2.

```
foldPal :: PalAlgebra p -> Pal -> p
foldPal (pal1, pal2, pal3, pal4, pal5) = fPal
 where
 fPal Pal1 = pal1
 fPal Pal2 = pal2
 fPal Pal3 = pal3
 fPal (Pal4 p) = pal4 (fPal p)
 fPal (Pal5 p) = pal5 (fPal p)
```

3.

```
a2cPal = foldPal ("
 , "a"
 , "b"
 , λp -> "a" ++ p ++ "a"
 , λp -> "b" ++ p ++ "b"
)
aCountPal = foldPal (0, 1, 0, λp -> p + 2, λp -> p)
```

4.

```
pfoldPal :: PalAlgebra p -> Parser Char p
pfoldPal (pal1, pal2, pal3, pal4, pal5) = pPal
 where
 pPal = const pal1 <$> ε
 <|> const pal2 <$> syma
 <|> const pal3 <$> symb
 <|> (_p -> pal4 p) <$> syma <*> pPal <*> syma
 <|> (_p -> pal5 p) <$> symb <*> pPal <*> symb
 syma = symbol 'a'
 symb = symbol 'b'
```

5. The parser *pfoldPal* *m*<sub>1</sub> returns the concrete representation of a palindrome.  
The parser *pfoldPal* *m*<sub>2</sub> returns the number of *a*'s occurring in a palindrome.

5.10

1. type MirAlgebra m = (m,m->m,m->m)
2. foldMir :: MirAlgebra m -> Mir -> m  
foldMir (mir1,mir2,mir3) = fMir where  
fMir Mir1 = mir1  
fMir (Mir2 m) = mir2 (fMir m)  
fMir (Mir3 m) = mir3 (fMir m)
3. a2cMir = foldMir (" , \m->"a"++m++"a" , \m->"b"++m++"b")  
m2pMir = foldMir (Pal1,Pal4,Pal5)

4.

```

pfoldMir :: MirAlgebra m -> Parser Char m
pfoldMir (mir1, mir2, mir3) = pMir where
 pMir = const mir1 <$> epsilon
 <|> (_ m _ -> mir2 m) <$> syma <*> pMir <*> syma
 <|> (_ m _ -> mir3 m) <$> symb <*> pMir <*> symb
 syma = symbol 'a'
 symb = symbol 'b'

```
5. The parser *pfoldMir*  $m_1$  returns the concrete representation of a palindrome. The parser *pfoldMir*  $m_2$  returns the abstract representation of a palindrome.

### 5.11

1.

```

type ParityAlgebra p = (p,p->p,p->p,p->p)

```
2.

```

foldParity :: ParityAlgebra p -> Parity -> p
foldParity (empty,parity1,parityL0,parityR0) = fParity where
 fParity Empty = empty
 fParity (Parity1 p) = parity1 (fParity p)
 fParity (ParityL0 p) = parityL0 (fParity p)
 fParity (ParityR0 p) = parityR0 (fParity p)

```
3.

```

a2cParity = foldParity ("
 ,\x->"1"++x++"1"
 ,\x->"0"++x
 ,\x->x++"0"
)

```

### 5.12

1.

```

type BitListAlgebra bl z b = ((b->z->bl,b->bl),(bl->z->z,bl->z),(b,b))

```
2.

```

foldBitList :: BitListAlgebra bl z b -> BitList -> bl
foldBitList ((consb,singleb),(consbl,singlebl),(bit0,bit1)) = fBitList
 where
 fBitList (ConsB b z) = consb (fBit b) (fZ z)
 fBitList (SingleB b) = singleb (fBit b)
 fZ (ConsBL bl z) = consbl (fBitList bl) (fZ z)
 fZ (SingleBL bl) = singlebl (fBitList bl)
 fBit Bit0 = bit0
 fBit Bit1 = bit1

```
3.

```

a2cBitList = foldBitList (((++),id)
 ,((++),id)
 ,("0","1")
)

```
4.

```

pfoldBitList :: BitListAlgebra bl z b -> Parser Char bl
pfoldBitList ((consb,singleb),(consbl,singlebl),(bit0,bit1)) = pBitList
 where
 pBitList = consb <$> pBit <*> pZ
 <|> singleb <$> pBit

```

B. Answers to exercises

```

pZ = (_ bl z -> consbl bl z) <$> symbol ','
 <*> pBitList
 <*> pZ

 <|> singlebl <$> pBitList
pBit = const bit0 <$> symbol '0'
 <|> const bit1 <$> symbol '1'

```

5.13

1.

```

data Block = B1 Stat Rest
data Rest = R1 Stat Rest | Nix
data Stat = S1 Decl | S2 Use | S3 Nest
data Decl = Dx | Dy
data Use = UX | UY
data Nest = N1 Block

```

The abstract representation of  $x; (y; Y); X$  is

```

B1 (S1 Dx) (R1 stat2 rest2)
where
 stat2 = S3 (N1 block)
 block = B1 (S1 Dy) (R1 (S2 UY) Nix)
 rest2 = R1 (S2 UX) Nix

```

2.

```

type BlockAlgebra b r s d u n =
 (s → r → b
 , (s → r → r, r)
 , (d → s, u → s, n → s)
 , (d, d)
 , (u, u)
 , b → n
)

```

3.

```

foldBlock :: BlockAlgebra b r s d u n → Block → b
foldBlock (b1, (r1, nix), (s1, s2, s3), (dx, dy), (ux, uy), n1) =
 foldB
where
 foldB (B1 stat rest) = b1 (foldS stat) (foldR rest)
 foldR (R1 stat rest) = r1 (foldS stat) (foldR rest)

```



```

foldR Nix = nix
foldS (S1 decl) = s1 (foldD decl)
foldS (S2 use) = s2 (foldU use)
foldS (S3 nest) = s3 (foldN nest)
foldD Dx = dx
foldD Dy = dy
foldU UX = ux
foldU UY = uy
foldN (N1 block) = n1 (foldB block)

```

4.

```

a2cBlock :: Block -> String
a2cBlock = foldBlock a2cBlockAlgebra
a2cBlockAlgebra :: BlockAlgebra String String String
 String String String
a2cBlockAlgebra =
 (b1, (r1, nix), (s1, s2, s3), (dx, dy), (ux, uy), n1)
where
 b1 u v = u ++ v
 r1 u v = ";" ++ u ++ v
 nix = ""
 s1 u = u
 s2 u = u
 s3 u = u
 dx = "x"
 dy = "y"
 ux = "X"
 uy = "Y"
 n1 u = "(" ++ u ++ ")"

```

**7.1** The type `TreeAlgebra` and the function `foldTree` are defined in the `rep-min` problem.

```

1. height = foldTree heightAlgebra
heightAlgebra = (const 0, \l r -> (1 'max' r) + 1)
--frontAtLevel :: Tree -> Int -> [Int]
--frontAtLevel (Leaf i) h = if h == 0 then [i] else []
--frontAtLevel (Bin l r) h =
-- if h > 0 then frontAtLevel l (h-1) ++ frontAtLevel r (h-1)
-- else []
frontAtLevel = foldTree frontAtLevelAlgebra
frontAtLevelAlgebra =
 (\i h -> if h == 0 then [i] else []

```

B. Answers to exercises

```
, \f g h -> if h > 0 then f (h-1) ++ g (h-1) else []
)
```

2. a) Straightforward Solution: impossible to give a solution like rm.sol1.

```
heightAlgebra = (const 0, \l r -> (l 'max' r) + 1)
front t = foldTree frontAtLevelAlgebra t h
 where h = foldTree heightAlgebra t
```

- b) Lambda Lifting

```
front' t = foldTree
 frontAtLevelAlgebra
 t
 (foldTree heightAlgebra t)
```

- c) Tupling Computations

```
htupfr :: TreeAlgebra (Int, Int -> [Int])
htupfr
-- = (heightAlgebra 'tuple' frontAtLevelAlgebra)
 = (\i -> (0
 , \h -> if h == 0 then [i] else []
)
 , \(\lh,f) (rh,g) -> ((lh 'max' rh) + 1
 , \h -> if h > 0
 then f (h-1) ++ g (h-1)
 else []
)
)
front'' t = fr h
 where (h, fr) = foldTree htupfr t
```

- d) Merging Tupled Functions

It is helpful to note that the merged algebra is constructed such that

```
foldTree mergedAlgebra t i = (height t, frontAtLevel t i)
```

Therefore, the definitions corresponding to the fourth solution are

```
mergedAlgebra :: TreeAlgebra (Int -> (Int, [Int]))
mergedAlgebra =
 (\i -> \h -> (0
 , if h == 0 then [i] else []
)
 , \lfun rfun -> \h -> let (lh, xs) = lfun (h-1)
 (rh, ys) = rfun (h-1)
 in
 ((lh 'max' rh) + 1
```

```

 , if h > 0 then xs ++ ys else []
)
)
 front''' t = fr
 where (h, fr) = foldTree mergedAlgebra t h

```

**7.2** The `highest_front` problem is the problem of finding the first non-empty list of nodes which are at the lowest level. The solutions are similar to these of the `deepest_front` problem.

```

1. --lowest :: Tree -> Int
 --lowest (Leaf i) = 0
 --lowest (Bin l r) = ((lowest l) 'min' (lowest r)) + 1
 lowest = foldTree lowAlgebra
 lowAlgebra = (const 0, \ l r -> (l 'min' r) + 1)

```

2. a) Straightforward Solution: impossible to give a solution like `rm.soll`.

```

lfront t = foldTree frontAtLevelAlgebra t l
 where l = foldTree lowAlgebra t

```

b) Lambda Lifting

```

lfront' t = foldTree
 frontAtLevelAlgebra
 t
 (foldTree lowAlgebra t)

```

c) Tupling Computations

```

ltupfr :: TreeAlgebra (Int, Int -> [Int])
ltupfr =
 (\i -> (0
 , (\l -> if l == 0 then [i] else [])
)
 , \ (lh,f) (rh,g) -> ((lh 'min' rh) + 1
 , \l -> if l > 0
 then f (l-1) ++ g (l-1)
 else []
)
)
lfront'' t = fr l
 where (l, fr) = foldTree ltupfr t

```

d) Merging Tupled Functions

It is helpful to note that the merged algebra is constructed such that

```

foldTree mergedAlgebra t i = (lowest t, frontAtLevel t i)

```

B. Answers to exercises

Therefore, the definitions corresponding to the fourth solution are

```
lmergedAlgebra :: TreeAlgebra (Int -> (Int, [Int]))
lmergedAlgebra =
 (\i -> \l -> (0
 , if l == 0 then [i] else []
)
 , \lfun rfun -> \l ->
 let (ll, lres) = lfun (l-1)
 (rl, rres) = rfun (l-1)
 in
 ((ll 'min' rl) + 1
 , if l > 0 then lres ++ rres else []
)
)
lfront''' t = fr
 where (l, fr) = foldTree lmergedAlgebra t l
```

**8.1** Let  $G = (T, N, R, S)$ . Define the regular grammar  $G' = (T, N \cup \{S'\}, R', S')$  as follows.  $S'$  is a new nonterminal symbol. For the productions  $R'$ , divide the productions of  $R$  in two sets: the productions  $R1$  of which the right-hand side consists just of terminal symbols, and the productions  $R2$  of which the right-hand side ends with a nonterminal. Define a set  $R3$  of productions by adding the nonterminal  $S$  at the right end of each production in  $R1$ . Define  $R' = R \cup R3 \cup S' \rightarrow S \mid \epsilon$ . The grammar  $G'$  thus defined is regular, and generates the language  $L^*$ .

**8.2** In step 2, add the productions

```
S → aS
S → ε
S → cC
S → a
A → ε
A → cC
A → a
B → cC
B → a
```

and remove the productions

```
S → A
A → B
B → C
```

In step 3, remove the productions  $A \rightarrow \epsilon, B \rightarrow \epsilon$ .

8.3

$$\begin{aligned} \text{ndfsa } d \text{ } qs [ ] &= qs \\ \text{ndfsa } d \text{ } qs (x : xs) &= \text{ndfsa } d (d \text{ } qs \ x) \ xs \end{aligned}$$

8.4 Let  $M = (X, Q, d, S, F)$  be a deterministic finite-state automaton. Then the nondeterministic finite-state automaton  $M'$  defined by  $M' = (X, Q, d', \{S\}, F)$ , where  $d' \ q \ x = \{d \ q \ x\}$  accepts the same language.

8.5 Let  $M = (X, Q, d, S, F)$  be a deterministic finite-state automaton that accepts language  $L$ . Then the deterministic finite-state automaton  $M' = (X, Q, d, S, Q - F)$ , where  $Q - F$  is the set of states  $Q$  from which all states in  $F$  have been removed, accepts the language  $\bar{L}$ . Here we assume that  $d \ q \ a$  is defined for all  $q$  and  $a$ .

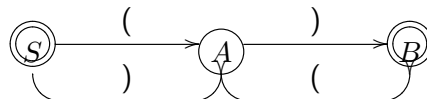
8.6

1.  $L_1 \cap L_2 = \overline{(\bar{L}_1 \cup \bar{L}_2)}$ , so it follows that regular languages are closed under intersection if they are closed under complementation and union. Regular languages are closed under union, see Theorem 8.10, and they are closed under complementation, see Exercise 8.5.

$$\begin{aligned} 2. \quad & \text{Ldfa } M \\ &= \\ & \{w \in X^* \mid \text{dfa\_accept } w (d, (S_1, S_2), (F_1 \times F_2))\} \\ &= \\ & \{w \in X^* \mid \text{dfa } d (S_1, S_2) \ w \in F_1 \times F_2\} \\ &= \quad \text{This requires a proof by induction} \\ & \{w \in X^* \mid (\text{dfa } d_1 \ S_1 \ w, \text{dfa } d_2 \ S_2 \ w) \in F_1 \times F_2\} \\ &= \\ & \{w \in X^* \mid \text{dfa } d_1 \ S_1 \ w \in F_1 \wedge \text{dfa } d_2 \ S_2 \ w \in F_2\} \\ &= \\ & \{w \in X^* \mid \text{dfa } d_1 \ S_1 \ w \in F_1\} \cap \{w \in X^* \mid \text{dfa } d_2 \ S_2 \ w \in F_2\} \\ &= \\ & \text{Ldfa } M_1 \cap \text{Ldfa } M_2 \end{aligned}$$

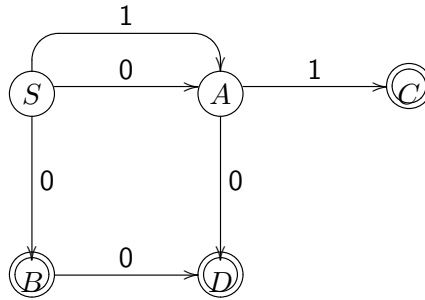
8.7

1.



B. Answers to exercises

2.



**8.8** Use the definition of  $Lre$  to calculate the languages:

1.  $\{\epsilon, b\}$
2.  $(bc^*)$
3.  $\{a\}b^* \cup c^*$

**8.9**

$$\begin{aligned}
 1. \quad & Lre(R(S + T)) \\
 &= \\
 & (Lre(R))(Lre(S + T)) \\
 &= \\
 & (Lre(R))(Lre(S) \cup Lre(T)) \\
 &= \\
 & (Lre(R))(Lre(S)) \cup (Lre(R))(Lre(T)) \\
 &= \\
 & Lre(RS + RT)
 \end{aligned}$$

2. Similar to the above calculation.

**8.10** Take  $R = a$ ,  $S = a$ , and  $R = a$ ,  $S = b$ .

**8.11** If both  $V$  and  $W$  are subsets of  $S$ , then  $Lre(R(S + V)) = Lre(R(S + W))$ . Since  $S \neq \emptyset$ ,  $V = S$  and  $W = \emptyset$  satisfy the requirement. Another solution is

$$\begin{aligned}
 V &= S \cap R \\
 W &= S \cap \overline{R}
 \end{aligned}$$

Since  $S \neq \emptyset$ , at least one of  $S \cap R$  and  $S \cap \overline{R}$  is not empty, and it follows that  $V \neq W$ . There exist other solutions than these.

**8.12** The string  $01$  may not occur in a string in the language of the regular expression. So when a  $0$  appears somewhere, only  $0$ 's can follow. Take  $1*0^*$ .

**8.13**

1. If we can give a regular grammar for  $(a + bb)^* + c$ , we can use the procedure constructed in Exercise 8.1 to obtain a regular grammar for the complete regular expression. The following regular grammar generates  $(a + bb)^* + c$ :

$$S_0 \rightarrow S_1 \mid c$$

provided  $S_1$  generates  $(a + bb)^*$ . Again, we use Exercise 8.1 and a regular grammar for  $a + bb$  to obtain a regular grammar for  $(a + bb)^*$ . The language of the regular expression  $a + bb$  is generated by the regular grammar

$$S_2 \rightarrow a \mid bb$$

2. The language of the regular expression  $a^*$  is generated by the regular grammar

$$S_0 \rightarrow \epsilon \mid aS_0$$

The language of the regular expression  $b^*$  is generated by the regular grammar

$$S_1 \rightarrow \epsilon \mid bS_1$$

The language of the regular expression  $ab$  is generated by the regular grammar

$$S_2 \rightarrow ab$$

It follows that the language of the regular expression  $a^* + b^* + ab$  is generated by the regular grammar

$$S_3 \rightarrow S_0 \mid S_1 \mid S_2$$

**8.14** First construct a nondeterministic finite-state automaton for these grammars, and use the automata to construct the following regular expressions.

1.  $a(b + b(b + ab)^*(ab + \epsilon)) + bb^* + \epsilon$
2.  $(0 + 10^*1)^*$

### 8.16

1.  $b^* + b^*(aa^*b(a + b)^*)$
2.  $(b + a(bb)^*b)^*$

**9.1** We follow the proof pattern for non-regular languages given in the text.

Let  $n \in \mathbb{N}$ .

Take  $s = a^{n^2}$  with  $x = \epsilon$ ,  $y = a^n$ , and  $z = a^{n^2-n}$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = a^p$ ,  $v = a^q$  and  $w = a^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 2$ , then

B. Answers to exercises

$$\begin{aligned}
 & xuv^2wz \notin L \\
 \Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
 & a^{p+2q+r}a^{n^2-n} \notin L \\
 \Leftarrow & p + q + r = n \text{ and } q > 0 \\
 & n^2 + q \text{ is not a square} \\
 \Leftarrow & \\
 & n^2 < n^2 + q < (n + 1)^2 \\
 \Leftarrow & \\
 & q < 2n + 1
 \end{aligned}$$

**9.2** Using the proof pattern.

Let  $n \in \mathbb{N}$ .

Take  $s = a^n b^{n+1}$  with  $x = a^n$ ,  $y = b^n$ , and  $z = b$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = b^p$ ,  $v = b^q$  and  $w = b^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 0$ , then

$$\begin{aligned}
 & xuwz \notin L \\
 \Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
 & a^n b^{p+r+1} \notin L \\
 \Leftarrow & \\
 & p + q + r \geq p + r + 1 \\
 \Leftarrow & \\
 & q > 0
 \end{aligned}$$

**9.3** Using the proof pattern.

Let  $n \in \mathbb{N}$ .

Take  $s = a^n b^{2n}$  with  $x = a^n$ ,  $y = b^n$ , and  $z = b^n$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = b^p$ ,  $v = b^q$  and  $w = b^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 2$ , then

$$\begin{aligned}
 & xuv^2wz \notin L \\
 \Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
 & a^n b^{2n+q} \notin L \\
 \Leftarrow & \\
 & 2n + q > 2n \\
 \Leftarrow & \\
 & q > 0
 \end{aligned}$$



**9.4** Using the proof pattern.

Let  $n \in \mathbb{N}$ .

Take  $s = a^6 b^{n+4} a^n$  with  $x = a^6 b^{n+4}$ ,  $y = a^n$ , and  $z = \epsilon$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = a^p$ ,  $v = a^q$  and  $w = a^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 6$ , then

$$\begin{aligned} & xuv^6wz \notin L \\ \Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\ & a^6 b^{n+4} a^{n+5q} \notin L \\ \Leftarrow & \\ & n + 5q > n + 4 \\ \Leftarrow & \\ & q > 0 \end{aligned}$$

**9.5** The length of a substring with  $a$ 's,  $b$ 's, and  $c$ 's is at least  $r + 2$ , and  $|vwx| \leq d \leq r$ .

**9.6** We follow the proof pattern for proving that a language is not context-free given in the text.

Let  $c, d \in \mathbb{N}$ .

Take  $z = a^{k^2}$  with  $k = \max(c, d)$ .

Let  $u, v, w, x, y$  be such that  $z = uvwxy$ ,  $|vx| > 0$  and  $|vwx| \leq d$

That is  $u = a^p$ ,  $v = a^q$ ,  $w = a^r$ ,  $x = a^s$  and  $y = a^t$  with  $p + q + r + s + t = k^2$ ,  $q + r + s \leq d$  and  $q + s > 0$ .

Take  $i = 2$ , then

$$\begin{aligned} & uv^2wx^2y \notin L \\ \Leftarrow & \text{ defn. } u, v, w, x, y, \text{ calculus} \\ & a^{k^2+q+s} \notin L \\ \Leftarrow & \\ & q + s > 0 \\ & k^2 + q + s \text{ is not a square} \\ \Leftarrow & \\ & k^2 < k^2 + q + s < (k + 1)^2 \\ \Leftarrow & \\ & q + s < 2k + 1 \\ \Leftarrow & \text{ defn. } k \\ & q + s \leq d \end{aligned}$$

**9.7** Using the proof pattern.

Let  $c, d \in \mathbb{N}$ .

B. Answers to exercises

Take  $z = a^k$  with  $k$  is prime and  $k > \max(c, d)$ .

Let  $u, v, w, x, y$  be such that  $z = uvwxy$ ,  $|vx| > 0$  and  $|vwx| \leq d$

That is  $u = a^p$ ,  $v = a^q$ ,  $w = a^r$ ,  $x = a^s$  and  $y = a^t$  with  $p + q + r + s + t = k$ ,  $q + r + s \leq d$  and  $q + s > 0$ .

Take  $i = k + 1$ , then

$$\begin{aligned} & wv^{k+1}wx^{k+1}y \notin L \\ \Leftrightarrow & \text{ defn. } u, v, w, x, y, \text{ calculus} \\ & a^{k+kq+ks} \notin L \\ \Leftrightarrow & \\ & k(1 + q + s) \text{ is not a prime} \\ \Leftrightarrow & \\ & q + s > 0 \end{aligned}$$

**9.8** Using the proof pattern.

Let  $c, d \in \mathbb{N}$ .

Take  $z = a^k b^k a^k b^k$  with  $k = \max(c, d)$ .

Let  $u, v, w, x, y$  be such that  $z = uvwxy$ ,  $|vx| > 0$  and  $|vwx| \leq d$

Note that our choice for  $k$  guarantees that substring  $vwx$  has one of the following shapes:

- $vwx$  consists of just a's, or just b's.
- $vwx$  contains both a's and b's.

Take  $i = 0$ , then

- If  $vwx$  consists of just a's, or just b's, then it is impossible to write the string  $uwv$  as  $ww$  for some string  $w$ , since only the number of terminals of one kind is decreased.
- If  $vwx$  contains both a's and b's, it lies somewhere on the border between a's and b's, or on the border between b's and a's. Then the string  $uwv$  can be written as

$$uwv = a^s b^t a^p b^q$$

for some  $s, t, p, q$ , respectively. At least one of  $s, t, p$  and  $q$  is less than  $k$ , while two of them are equal to  $k$ . Again this sentence is not an element of the language.

**9.9** Using the proof pattern.

Let  $n \in \mathbb{N}$ .

Take  $s = 1^n 0^n$  with  $x = 1^n$ ,  $y = 0^n$ , and  $z = \epsilon$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = b^p$ ,  $v = b^q$  and  $w = b^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 2$ , then

$$\begin{aligned}
& xv^2w \notin L \\
\Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
& 1^n 0^{p+2q+r} \notin L \\
\Leftarrow & p + q + r = n \\
& 1^n 0^{n+q} \notin L \\
\Leftarrow & q > 0 \\
& \text{true}
\end{aligned}$$

### 9.10

- The language  $\{a^i b^j \mid 0 \leq i \leq j\}$  is context-free. The language can be generated by the following context-free grammar:

$$\begin{aligned}
S & \rightarrow AB \\
A & \rightarrow \epsilon \mid aAb \\
B & \rightarrow \epsilon \mid bB
\end{aligned}$$

This grammar generates the strings  $a^m b^m b^n$  for  $m, n \in \mathbb{N}$ . These are exactly the strings of the given language.

- The language is not regular. We prove this using the proof pattern.

Let  $n \in \mathbb{N}$ .

Take  $s = a^n b^{n+1}$  with  $x = \epsilon$ ,  $y = a^n$ , and  $z = b^{n+1}$ .

Let  $u, v, w$  be such that  $y = uvw$  with  $v \neq \epsilon$ , that is,  $u = a^p$ ,  $v = a^q$  and  $w = a^r$  with  $p + q + r = n$  and  $q > 0$ .

Take  $i = 3$ , then

$$\begin{aligned}
& xuv^3wz \notin L \\
\Leftarrow & \text{ defn. } x, u, v, w, z, \text{ calculus} \\
& a^{p+3q+r} b^{n+1} \notin L \\
\Leftarrow & p + q + r = n \\
& n + 2q > n + 1 \\
\Leftarrow & q > 0 \\
& \text{true}
\end{aligned}$$

### 9.11

- The language  $\{wcw \mid w \in \{a, b\}^*\}$  is not context-free. We prove this using the proof pattern.

Let  $c, d \in \mathbb{N}$ .

## B. Answers to exercises

Take  $z = a^k b^k c a^k b^k$  with  $k = \max(c, d)$ .

Let  $u, v, w, x, y$  be such that  $z = uvwxy$ ,  $|vx| > 0$  and  $|vwx| \leq d$

Note that our choice for  $k$  guarantees that  $|vwx| \leq k$ . The substring  $vwx$  has one of the following shapes:

- $vwx$  does not contain  $c$ .
- $vwx$  contains  $c$ .

Take  $i = 0$ , then

- If  $vwx$  does not contain  $c$  it is a substring at the left-hand side or at the right-hand side of  $c$ . Then it is impossible to write the string  $uvw$  as  $scs$  for some string  $s$ , since only the number of terminals on one side of  $c$  is decreased.
- If  $vwx$  contains  $c$  then the string  $vwx$  can be written as

$$vwx = b^s c a^t$$

for some  $s, t$  respectively. For  $uvw$  there are two possibilities.

The string  $uvw$  does not contain  $c$  so it is not an element of the language.

If the string  $uvw$  does contain  $c$  then it has either fewer  $b$ 's at the left-hand side than at the right-hand side or fewer  $a$ 's at the right-hand side than at the left-hand side, or both. So it is impossible to write the string  $uvw$  as  $scs$  for some string  $s$ .

2. The language is not regular because it is not context-free. The set of regular languages is a subset of the set of context-free languages.

### 9.12

1. The grammar  $G$  is context-free because there is only one nonterminal at the left hand side of the production rules and the right hand side is a sequence of terminals and nonterminals. The grammar is not regular because there is a production rule with two nonterminals at the right hand side.
2.  $L(G) = \{(st)^* \mid s \in \{^* \wedge t \in \}^* \wedge |s| = |t|\}$   
 $L(G)$  is all sequences of nested braces.
3. The language is context-free because it is generated by a context-free grammar. The language is not regular. We use the proof pattern and choose  $s = \{^n\}^n$ . The proof is exactly the same as the proof for non-regularity of  $L = \{a^m b^m \mid m \geq 0\}$  in Section 9.2.

### 9.13

1. The grammar is context-free. The grammar is not right-regular but left-regular because the nonterminal appears at the left hand side in the last production rule.

2.  $L(G) = 0^* \cup 10^*$

Note that the language is also generated by the following right-regular grammar:

$$\begin{aligned} S &\rightarrow \epsilon \\ S &\rightarrow 1A \\ S &\rightarrow 0A \\ A &\rightarrow \epsilon \\ A &\rightarrow 0 \\ A &\rightarrow 0A \end{aligned}$$

3. The language is context-free and regular because it is generated by a regular grammar.

**10.1** For both grammars we have:

```
empty = const False
```

**10.2** For grammar1 we have:

```
firsts S = {b,c}
firsts A = {a,b,c}
firsts B = {a,b,c}
firsts C = {a,b}
```

For grammar2:

```
firsts S = {a}
firsts A = {b}
```

**10.3** For gramm1 we have:

```
follow S = {a,b,c}
follow A = {a,b,c}
follow B = {a,b}
follow C = {a,b,c}
```

For gramm2:

```
follow = const {}
```

**10.4** For the productions of gramm3 we have:

B. Answers to exercises

```
lookAhead (S → AaS) = {a, c }
lookAhead (S → B) = {b}
lookAhead (A → cS) = {c}
lookAhead (A → []) = {a}
lookAhead (B → b) = {b}
```

Since all `lookAhead` sets for productions of the same nonterminal are disjoint, `gramm3` is an LL(1) grammar.

**10.5** After left factoring `gramm2` we obtain `gramm2'` with productions

```
S → aC
C → bA | a
A → bD
D → b | S
```

For this transformed grammar `gramm2'` we have:

The empty function

```
empty = const False
```

The first function

```
firsts S = {a}
firsts C = {a,b}
firsts A = {b}
firsts D = {a,b}
```

The follow function

```
follow = const {}
```

The `lookAhead` function

```
lookAhead (S → aC) = {a}
lookAhead (C → bA) = {b}
lookAhead (C → a) = {a}
lookAhead (A → bD) = {b}
lookAhead (D → b) = {b}
lookAhead (D → S) = {a}
```

Clearly, `gramm2'` is an LL(1) grammar.

**10.6** For the `empty` function we have:

```
empty R = True
empty _ = False
```

For the `firsts` function we have

```
firsts L = {0,1}
firsts R = {,}
firsts B = {0,1}
```

The follow function is

```
follow B = {,}
follow _ = {}
```

The `lookAhead` function is

```
lookAhead (L → B R) = {0, 1}
lookAhead (R → ε) = {}
lookAhead (R → , B R) = {,}
lookAhead (B → 0) = {0}
lookAhead (B → 1) = {1}
```

Since all `lookAhead` sets for productions of the same nonterminal are disjoint, the grammar is LL(1).

### 10.7

```
Node S [Node c []
 , Node A [Node c []
 , Node B [Node c [], Node c []]
 , Node C [Node b [], Node a []]
]
]
```

### 10.8

```
Node S [Node a []
 , Node C [Node b []
 , Node A [Node b [], Node D [Node b []]]
]
]
```

### 10.9

B. Answers to exercises

```
Node S [Node A []
 , Node a []
 , Node S [Node A [Node c [], Node S [Node B [Node b []]]]
 , Node a []
 , Node S [Node B [Node b []]]
]
]
```

10.10

1. 

```
list2Set :: Ord s => [s] -> [s]
list2Set = unions . map single
```
2. 

```
list2Set :: Ord s => [s] -> [s]
list2Set = foldr op []
 where
 op x xs = single x 'union' xs
```
3. 

```
pref :: Ord s => (s -> Bool) -> [s] -> [s]
pref p = list2Set . takeWhile p
```

or

```
pref p [] = []
pref p (x:xs) = if p x then single x 'union' pref p xs else []
```

or

```
pref p = foldr op []
 where
 op x xs = if p x then single x 'union' xs else []
```
4. 

```
prefplus p [] = []
prefplus p (x:xs) = if p x then single x 'union' prefplus p xs
 else single x
```
5. 

```
prefplus p = foldr op []
 where
 op x us = if p x then single x 'union' us
 else single x
```
6. 

```
prefplus p
=
 foldr op []
 where
 op x us = if p x then single x 'union' us else single x
=
 foldr op []
```



```

where
op x us = single x 'union' rest
 where
rest = if p x then us else []
=
 foldrRhs p single []

```

7. The function `foldrRhs p f []` takes a list `xs` and returns the set of all `f`-images of the elements of the prefix of `xs` all of whose elements satisfy `p` together with the `f`-image of the first element of `xs` that does not satisfy `p` (if this element exists).

The function `foldrRhs p f start` takes a list `xs` and returns the set `start 'union' foldrRhs p f [] xs`.

8.

## 10.11

1. For `gramm1`

```

a) foldrRhs empty first [] bSA
=
 unions (map first (prefplus empty bSA))
=
 empty b = False
 unions (map first [b])
=
 unions [{b}]
=
 {b}

```

```

b) foldrRhs empty first [] Cb
=
 unions (map first (prefplus empty Cb))
=
 empty C = False
 unions (map first [C])
=
 unions [{a, b}]
=
 { a, b }

```

2. For `gramm3`

```

a) foldrRhs empty first [] AaS
=
 unions (map first (prefplus empty AaS))

```

## B. Answers to exercises

```
= empty A = True
 unions (map first ([A] 'union' prefplus empty aS))
 empty a = False
 unions (map first ([A] 'union' [a]))
=
 unions (map first ([A, a]))
=
 unions [first A, first a]
=
 unions [{c}, {a}]
=
 {a, c}

b) scanRrhs empty first [] AaS
=
 map (foldrRhs empty first []) (tails AaS)
=
 map (foldrRhs empty first []) [AaS, aS, S, []]
=
 [foldrRhs empty first [] AaS
 , foldrRhs empty first [] aS
 , foldrRhs empty first [] S
 , foldrRhs empty first [] []
]
=
 calculus
 [{a, c}, {a}, {a, b, c }, []]
```

**11.1** Recall, for `gramm1` we have:

```
follow S = {a,b,c}
follow A = {a,b,c}
follow B = {a,b}
follow C = {a,b,c}
```

The production rule  $B \rightarrow cc$  cannot be used for a reduction when the input begins with the symbol  $c$ .

| stack | input  |
|-------|--------|
|       | ccccba |
| c     | cccba  |
| cc    | ccba   |
| ccc   | cba    |
| cccc  | ba     |
| Bcc   | ba     |
| bBcc  | a      |
| abBcc |        |
| CBcc  |        |
| Ac    |        |
| S     |        |

**11.2** Recall, for gramm2 we have:

follow S = {}  
follow A = {}

No reduction can be made until the input is empty.

| stack | input |
|-------|-------|
|       | abbb  |
| a     | bbb   |
| ba    | bb    |
| bba   | b     |
| bbba  |       |
| Aba   |       |
| S     |       |

**11.3** Recall, for gramm3 we have:

follow S = {a}  
follow A = {a}  
follow B = {a}

At the beginning a reduction is applied using the production rule  $A \rightarrow \varepsilon$ .

B. Answers to exercises

| stack | input |
|-------|-------|
|       | acbab |
| A     | acbab |
| aA    | cbab  |
| caA   | bab   |
| bcaA  | ab    |
| BcaA  | ab    |
| ScaA  | ab    |
| AaA   | ab    |
| aAaA  | b     |
| baAaA |       |
| SaAaA |       |
| SaA   |       |
| S     |       |