

Introduction to Functional Programming and the Structure of Programming Languages using OCaml

Lecture Notes
Version of October 25, 2023

Gert Smolka
Saarland University

Copyright © 2021 by Gert Smolka, all rights reserved

Contents

Preface	v
1 Getting Started	1
1.1 Programs and Declarations	1
1.2 Functions and Let Expressions	3
1.3 Conditionals, Comparisons, and Booleans	4
1.4 Recursive Power Function	5
1.5 Integer Division	6
1.6 Mathematical Level versus Coding Level	9
1.7 More about Mathematical Functions	11
1.8 A Higher-Order Function for Linear Search	13
1.9 Partial Applications	15
1.10 Inversion of Strictly Increasing Functions	17
1.11 Tail Recursion	18
1.12 Tuples	20
1.13 Exceptions and Spurious Arguments	21
1.14 Summary	23
2 Polymorphic Functions and Iteration	24
2.1 Polymorphic Functions	24
2.2 Iteration	25
2.3 Iteration on Pairs	26
2.4 Computing Primes	28
2.5 Polymorphic Typing Rules	30
2.6 Polymorphic Exception Raising and Equality Testing	32
2.7 Summary	33
3 Lists	35
3.1 Nil and Cons	35
3.2 Basic List Functions	36
3.3 List Functions in OCaml	39
3.4 Fine Points About Lists	41
3.5 Membership and List Quantification	42
3.6 Head and Tail	43
3.7 Position Lookup	44
3.8 Option Types	45
3.9 Generalized Match Expressions	46

Contents

3.10 Sublists	48
3.11 Folding Lists	49
3.12 Insertion Sort	51
3.13 Generalized Insertion Sort	53
3.14 Lexical Order	54
3.15 Prime Factorization	55
3.16 Key-Value Maps	57

Preface

This text teaches functional programming and the structure of programming languages to beginning students. It is written for the Programming 1 course for computer science students at Saarland University. We assume that incoming students are familiar with mathematical thinking, but we do not assume programming experience. The course is designed to take about one third of the study time of the first semester.

We have been teaching a course like this at Saarland University since 1998. Students perceive the course as challenging and exciting, whether they have programmed before or not. In 2021, we changed the teaching language to English and the programming language to OCaml.

As it comes to functional programming, we cover higher-order recursive functions, polymorphic typing, and constructor types for lists, trees, and abstract syntax. We emphasize the role of correctness statements and practice inductive correctness proofs. We also cover asymptotic running time considering binary search (logarithmic), insertion sort (quadratic), merge sort (linearithmic), and other algorithms.

As it comes to the structure of programming languages, we study the different layers of syntax and semantics at the example of the idealized functional programming language Mini-OCaml. We describe the syntactic layers with grammars and the semantic layers with inference rules. Based on these formal descriptions, we program recursive descent parsers, type checkers and evaluators.

We also cover stateful programming with arrays and cells (assignable variables). We explain how lists and trees can be stored as linked blocks in an array, thus explaining memory consumption for constructor types.

There is a textbook¹ written for the German iterations of the course (1998 to 2020). The new English text realizes some substantial changes: OCaml rather than Standard ML as programming language, less details about the concrete programming language being used, more emphasis on correctness arguments and algorithms, informal type-theoretic explanations rather than formal set-theoretic definitions.

The current version of the text leaves room for improvement. More basic explanations with more examples could be helpful in many places. An additional chapter on imperative programming with loops and the

¹Gert Smolka, *Programmierung — eine Einführung in die Informatik mit Standard ML*. Oldenbourg Verlag, 2008.

Preface

realization with stack machines with jumps (see the German textbook) would be interesting, but this extra material may not fit into the time-budget of a one-semester course.

At Saarland University, the course spans 15 weeks of teaching in the winter semester. Each week comes with two lectures (90 minutes each), an exercises assignment, office hours, tutorials, and a test (15 minutes). There is a midterm exam in December and a final exam (offered twice) after the lecture period has finished. The 2021 iteration of the course also came with a take home project over the holiday break asking the students to write an interpreter for Mini-OCaml. The take home project should be considered an important part of the course, given that it requires the students writing and debugging a larger program (about 250 lines), which is quite different from the small programs (up to 10 lines) the weekly assignments ask for.

1 Getting Started

In this chapter we start programming in OCaml. We use an interactive tool called an interpreter that checks and executes the declarations of a program one by one. We concentrate on recursive functions on integers computing things like powers, integer quotients, digit sums, and integer roots. We also formulate a general algorithm known as linear search as a higher-order function. We follow an approach known as functional programming where functions are designed at a mathematical level using types and equations before they are coded in a concrete programming language.

1.1 Programs and Declarations

An OCaml **program** is a sequence of **declarations**, which are executed in the order they are written. Our first program

```
let a = 2 * 3 + 2
let b = 5 * a
```

consists of two declarations. The first declaration **binds** the **identifier** *a* to the integer 8, and the second declaration binds the identifier *b* to the integer 40. This is clear from our understanding of the arithmetic operations “+” and “*”.

To learn programming in OCaml, you want to use an interactive tool called an **interpreter**.¹ The user feeds the interpreter with text. The interpreter checks that the given text can be interpreted as a program formulated correctly according to the rules of the programming language. If this is the case, the interpreter determines a **type** for every identifier and every expression of the program. Our example program is formulated correctly and the declared identifiers *a* and *b* both receive the type **int** (for integer). After a program has been checked successfully, the interpreter will execute it. In our case, execution will bind the identifier *a* to the integer 8 and the identifier *b* to the integer 40. After the program has been executed successfully, the interpreter will show the values it has computed for the declared identifiers. If a program is not formulated correctly, the interpreter will show an error message indicating which rule of the language is violated. Once the interpreter

¹A nice browser-based interpreter is <https://try.ocamlpro.com>.

1 Getting Started

has checked and executed a program, the user can extend it with further declarations. This way one can write the declarations of a program one by one.

At this point you want to start working with the interpreter. You will learn the exact rules of the language through experiments with the interpreter guided by the examples and explanations given in this text.

Here is a program redeclaring the identifier *a*:

```
let a = 2 * 3 + 2
let b = 5 * a
let a = 5
let c = a + b
```

The second declaration of the identifier *a* **shadows** the first declaration of *a*. Shadowing does not affect the binding of *b* since it is obtained before the second declaration of *a* is executed. After execution of the program, the identifiers *a*, *b*, and *c* are bound to the integers 5, 40, and 45, respectively.

The declarations we consider in this chapter all start with the **keyword** *let* and consist of a **head** and a **body** separated by the equality symbol “=”. Keywords cannot be used as identifiers. The bodies of declarations are **expressions**. Expressions can be obtained with identifiers, constants, and operators. The **nesting of expressions** can be arranged with **parentheses**. For instance, the expression $2 \cdot 3 + 2 - x$ may be written with **redundant parentheses** as $((2 \cdot 3) + 2) - x$. The parentheses in $2 \cdot (x + y) - 3$ are not redundant and are needed to make the expression $x + y$ the right argument of the product with the left argument 2.

Every expression has a **type**. So far we have only seen expressions of the type *int*. The values of the type *int* are integers (whole numbers $\dots, -2, -1, 0, 1, -2, \dots$). In contrast to the infinite mathematical type \mathbb{Z} , OCaml’s type *int* provides only a finite interval of **machine integers** realized efficiently by the hardware of the underlying computer. The endpoints of the interval can be obtained with the predefined identifiers *min_int* and *max_int*. All machine operations respect this interval. We have $\text{max_int} + 1 = \text{min_int}$, for instance.²

When we reason about programs, we will usually ignore the machine integers and just assume that all integers are available. As long as the

²It turns out that different interpreters realize different intervals for machine integers, even on the same computer. For instance, on the author’s computer, in October 2021, the browser-based Try OCaml interpreter realizes *max_int* as $2^{31} - 1 = 2147483647$, while the official OCaml interpreter realizes *max_int* as $2^{61} - 1 = 4611686018427387903$.

1 Getting Started

numbers in a concrete computation are small enough, this simplification does not lead to wrong conclusions.

Every programming language provides machine integers for efficiency. There are techniques for realizing much larger intervals of integers based on machine integers in programming languages.

Exercise 1.1.1 Give an expression and a declaration. Explain the structure of a declaration. Explain nesting of expressions. Give a type.

Exercise 1.1.2 To what value does the program

```
let a = 2 let a = a * a let a = a * a
```

bind the identifier a ?

Exercise 1.1.3 Give a machine integer x such that $x + 1 < x$.

1.2 Functions and Let Expressions

Things become interesting once we declare functions. The declaration

```
let square x = x * x
```

declares a function

$$\text{square} : \text{int} \rightarrow \text{int}$$

squaring its argument. The identifier *square* receives a **functional type** $\text{int} \rightarrow \text{int}$ describing functions that given an integer return an integer. Given the declaration of *square*, execution of the declaration

```
let a = square 5
```

binds the identifier a to the integer 25.

Can we compute a power x^8 with just three multiplications? Easy, we just square the integer x three times:

```
let pow8 x = square (square (square x))
```

This declaration gives us a function $\text{pow8} : \text{int} \rightarrow \text{int}$.

Another possibility is the declaration

```
let pow8' x =  
  let a = x * x in  
  let b = a * a in  
  b * b
```

1 Getting Started

declaring a function $pow8' : int \rightarrow int$ doing the three multiplications using two **local declarations**. Local declarations are obtained with **let expressions**

$let\ d\ in\ e$

combining a declaration d with an expression e using the keywords *let* and *in*. Note that the body of pow' nests two let expressions as $let\ d_1\ in\ (let\ d_2\ in\ e)$. OCaml lets you write redundant parentheses marking the nesting (if you want to). Let expressions must not be confused with top-level declarations (which don't use the keyword *in*).

Exercise 1.2.1 Write a function computing x^5 with just 3 multiplications. Write both a version with *square* and a version with local declarations. Write the version with local declarations with and without redundant parentheses marking the nesting.

1.3 Conditionals, Comparisons, and Booleans

The declaration

```
let abs x = if x < 0 then -x else x
```

declares a function $abs : int \rightarrow int$ returning the absolute value of an integer (e.g., $abs(-5) = 5$). The declaration of abs uses a **comparison** $x < 0$ and a **conditional** formed with the keywords **if**, **then**, and **else**.

The declaration

```
let max x y : int = if x <= y then y else x
```

declares a function $max : int \rightarrow int \rightarrow int$ computing the maximum of two numbers. This is the first function we see taking two arguments. There is an explicit **return type specification** in the declaration (appearing as “: *int*”), which is needed so that max receives the correct type.³

Given max , we can declare a function

```
let max3 x y z = max (max x y) z
```

returning the maximum of three numbers. This time no return type specification is needed since max forces the correct type $max3 : int \rightarrow int \rightarrow int \rightarrow int$.

Next we declare a three-argument maximum function using a local declaration:

³The complication stems from the design that in OCaml comparisons like \leq also apply to types other than *int*.

1 Getting Started

```
let max3 x y z : int =  
  let a = if x <= y then y else x in  
  if a <= z then z else a
```

There is a type *bool* having two values *true* and *false* called **booleans**. The comparison operator “ \leq ” used for *max* (written “ \leq ” in OCaml) is used with the functional type $int \rightarrow int \rightarrow bool$ saying that an expression $e_1 \leq e_2$ where both subexpressions e_1 and e_2 have type *int* has type *bool*.

The declaration

```
let test (x : int) y z = if x <= y then y <= z else false
```

declares a function $test : int \rightarrow int \rightarrow int \rightarrow bool$ testing whether its three arguments appear in order. The **type specification** given for the first argument x is needed so that *test* receives the correct type. Alternatively, type specifications could be given for one or all of the other arguments.

Exercise 1.3.1 Declare minimum functions analogous to the maximum functions declared above. For each declared function give its type (before you check it with the interpreter). Also, write the declarations with and without redundant parentheses to understand the nesting.

Exercise 1.3.2 Declare functions $int \rightarrow int \rightarrow bool$ providing the comparisons $x = y$, $x \neq y$, $x < y$, $x \leq y$, $x > y$, and $x \geq y$. Do this by just using conditionals and comparisons $x \leq y$. Then realize $x \leq y$ with $x \leq 0$ and subtraction.

1.4 Recursive Power Function

Our next goal is a function computing powers x^n using multiplication. We recall that powers satisfy the equations

$$\begin{aligned}x^0 &= 1 \\ x^{n+1} &= x \cdot x^n\end{aligned}$$

Using the equations, we can compute every power x^n where $n \geq 0$. For instance,

$$\begin{aligned}2^3 &= 2 \cdot 2^2 && \text{2nd equation} \\ &= 2 \cdot 2 \cdot 2^1 && \text{2nd equation} \\ &= 2 \cdot 2 \cdot 2 \cdot 2^0 && \text{2nd equation} \\ &= 2 \cdot 2 \cdot 2 \cdot 1 && \text{1st equation} \\ &= 8\end{aligned}$$

1 Getting Started

The trick is that the 2nd equation reduces larger powers to smaller powers, so that after repeated application of the 2nd equation the power x^0 appears, which can be computed with the 1st equation. What we see here is a typical example of a recursive computation.

Recursive computations can be captured with **recursive functions**. To arrive at a function computing powers, we merge the two equations for powers into a single equation using a conditional:

$$x^n = \text{IF } n < 1 \text{ THEN } 1 \text{ ELSE } x \cdot x^{n-1}$$

We now declare a recursive function implementing this equation:

```
let rec pow x n =  
  if n < 1 then 1  
  else x * pow x (n - 1)
```

The identifier *pow* receives the type $\text{pow} : \text{int} \rightarrow \text{int} \rightarrow \text{int}$. We say that the function *pow* applies itself. Recursive function applications are admitted if the declaration of the function uses the keyword **rec**.

We have demonstrated an important point about programming at the example of the power function: Recursive functions are designed at a mathematical level using equations. Once we have the right equations, we can implement the function in any programming language.

1.5 Integer Division

Given two integers $x \geq 0$ and $y > 0$, there exist unique integers $k, r \geq 0$ such that

$$x = k \cdot y + r$$

and

$$r < y$$

We call k the **quotient** and r the **remainder** of x and y . For instance, given $x = 11$ and $y = 3$, we have the quotient 3 and the remainder 2 (since $11 = 3 \cdot 3 + 2$). We speak of *integer division*, or *division with remainder*, or *Euclidean division*.

We may also characterize the quotient as the largest number k such that $k \cdot y \leq x$, and define the remainder as $r = x - k \cdot y$.

There is a nice geometrical interpretation of integer division. The idea is to place boxes of the same length y next to each other into a shelf of length x . The maximal number of boxes that can be placed into the shelf is the quotient k and the length of the space remaining is the

1 Getting Started

remainder $r = x - k \cdot y < y$. For instance, if the shelf has length 11 and each box has length 3, we can place at most 3 boxes into the shelf, with 2 units of length remaining.⁴

Given that x and y uniquely determine k and r , we are justified in using the notations x/y and $x \% y$ for k and r . By definition of k and r , we have

$$x = (x/y) \cdot y + x \% y$$

and

$$x \% y < y$$

for all $x \geq 0$ and $y > 0$.

From our explanations it is clear that we can compute x/y and $x \% y$ given x and y . In fact, the resulting operations x/y and $x \% y$ are essential for programming and are realized efficiently for machine integers on computers. We refer to the operations as *division* and *modulo*, or just as “div” and “mod”. Accordingly, we read the applications x/y and $x \% y$ as “ x div y ” and “ x mod y ”. OCaml provides both operations as primitive operations using the notations x/y and $x \bmod y$.

Digit sum

With div and mod we can decompose the decimal representation of numbers. For instance, $367 \% 10 = 7$ and $367/10 = 36$. More generally, $x \% 10$ yields the last digit of the decimal representation of x , and $x/10$ cuts off the last digit of the decimal representation of x .

Knowing these facts, we can declare a recursive function computing the digit sum of a number:

```
let rec digit_sum x =  
  if x < 10 then x  
  else digit_sum (x / 10) + (x mod 10)
```

For instance, we have $\text{digit_sum } 367 = 16$. We note that *digit_sum* terminates since the argument gets smaller upon recursion.

Exercise 1.5.1 (First digit) Declare a function that yields the first digit of the decimal representation of a number. For instance, the first digit of 367 is 3.

Exercise 1.5.2 (Maximal digit) Declare a function that yields the maximal digit of the decimal representation of a number. For instance, the maximal digit of 376 is 7.

⁴A maybe simpler geometrical interpretation of integer division asks how many boxes of height y can be stacked on each other without exceeding a given height x .

1 Getting Started

Digit reversal

We now write a function *rev* that given a number computes the number represented by the reversed digital representation of the number. For instance, we want $rev\ 76 = 67$, $rev\ 67 = 76$, and $rev\ 7600 = 67$. To write *rev*, we use an important algorithmic idea. The trick is to have an additional **accumulator argument** that is initially 0 and that collects the digits we cut off at the right of the main argument. For instance, we want the **trace**

$$rev'\ 456\ 0 = rev'\ 45\ 6 = rev'\ 4\ 65 = rev'\ 0\ 654 = 654$$

for the helper function *rev'* with the accumulator argument.

We declare *rev* and the helper function *rev'* as follows:

```
let rec rev' x a =  
  if x <= 0 then a  
  else rev' (x / 10) (10 * a + x mod 10)  
let rev x = rev' x 0
```

We refer to *rev'* as the **worker function** for *rev* and to the argument *a* of *rev'* as the **accumulator argument** of *rev'*. We note that *rev'* terminates since the first argument gets smaller upon recursion.

Greatest common divisors

Recall the notion of greatest common divisors. For instance, the greatest common divisor of 34 and 85 is the number 17. In general, two numbers $x, y \geq 0$ such that $x + y > 0$ always have a unique greatest common divisor. We assume the following rules for greatest common divisors (gcds for short):⁵

1. The gcd of x and 0 is x .
2. If $y > 0$, the gcd of x and y is the gcd of y and $x \% y$.

The two rules suffice to declare a function $gcd : int \rightarrow int \rightarrow int$ computing the gcd of two numbers $x, y \geq 0$ such that $x + y > 0$:

```
let rec gcd x y =  
  if y < 1 then x  
  else gcd y (x mod y)
```

The function terminates for valid arguments since $(x \% y) < y$ for $x \geq 0$ and $y \geq 1$.

Computing div and mod with repeated subtraction

We can compute x/y and $x \% y$ using repeated subtraction. To do so, we simply subtract y from x as long as we do not obtain a negative

⁵We will prove the correctness of the rules in a later chapter.

1 Getting Started

number. Then x/y is the number of successful subtractions and b is the remaining number.

```
let rec my_div x y = if x < y then 0 else 1 + my_div (x - y) y
let rec my_mod x y = if x < y then x else my_mod (x - y) y
```

We remark that both functions terminate for $x \geq 0$ and $y > 0$ since the first argument gets smaller upon recursion.

Exercise 1.5.3 (Traces) Give traces for the following applications:

```
rev' 678 0    rev' 6780 0    gcd 90 120    gcd 153 33    my_mod 17 5
```

We remark that the functions rev' , gcd , and my_mod employ a special form of recursion known as *tail recursion*. We will discuss tail recursion in §1.11.

1.6 Mathematical Level versus Coding Level

When we design a function for OCaml or another programming language, we do this at the mathematical level. The same is true when we reason about functions and their correctness. Designing and reasoning at the mathematical level has the advantage that it serves all programming languages, not just the concrete programming language we have chosen to work with (OCaml in our case). Given the design of a function at the mathematical level, we refer to the realization of the function in a concrete programming language as *coding*. Programming as we understand it emphasizes design and reasoning over coding.

It is important to distinguish between the mathematical level and the coding level. At the mathematical level, we ignore the type *int* of machine integers and instead work with infinite *mathematical types*. In particular, we will use the types

\mathbb{N} : 0, 1, 2, 3, ...	natural numbers
\mathbb{N}^+ : 1, 2, 3, 4, ...	positive integers
\mathbb{Z} : ..., -2, -1, 0, 1, 2, ...	integers
\mathbb{B} : false, true	booleans

When start with the design of a function, it is helpful to fix a mathematical type for the function. Once the type is settled, we can collect equations the function should satisfy. The goal here is to come up with a collection of equations that is sufficient for computing the function.

For instance, when we design a power function, we may start with the mathematical type

$$pow : \mathbb{Z} \rightarrow \mathbb{N} \rightarrow \mathbb{Z}$$

1 Getting Started

and the equation

$$\text{pow } x \ n = x^n$$

Together, the type and the equation *specify* the function we want to define. Next we need equations that can serve as defining equations for *pow*. The specifying equation is not good enough since we assume, for the purpose of the example, that the programming language we want to code in doesn't have a power operator. We now recall that powers x^n satisfy the equations

$$\begin{aligned} x^0 &= 1 \\ x^{n+1} &= x \cdot x^n \end{aligned}$$

In §1.4 we have already argued that rewriting with the two equations suffices to compute all powers we can express with the type given for *pow*. Next we adapt the equations to the function *pow* we are designing:

$$\begin{aligned} \text{pow } x \ 0 &= 1 \\ \text{pow } x \ n &= x \cdot \text{pow } x \ (n - 1) && \text{if } n > 0 \end{aligned}$$

The second equation now comes with an **application condition** replacing the pattern $n + 1$ in the equation $x^{n+1} = x \cdot x^n$.

We observe that the equations are **exhaustive** and **disjoint**, that is, for all x and n respecting the type of *pow*, the left side of one and only one of the equations applies to the **application** $\text{pow } x \ n$. We choose the equations as **defining equations** for *pow* and summarize our design with the mathematical definition

$$\begin{aligned} \text{pow} : \mathbb{Z} &\rightarrow \mathbb{N} \rightarrow \mathbb{Z} \\ \text{pow } x \ 0 &:= 1 \\ \text{pow } x \ n &:= x \cdot \text{pow } x \ (n - 1) \quad \text{if } n > 0 \end{aligned}$$

Note that we write defining equations with the symbol “:=” to mark them as defining.

We observe that the defining equations for *pow* are **terminating**. The **termination argument** is straightforward: Each **recursion step** issued by the second equation decreases the second argument n by 1. This ensures termination since a chain $x_1 > x_2 > x_3 > \dots$ of natural numbers cannot be infinite.

Next we code the mathematical definition as a function declaration in OCaml:

1 Getting Started

```
let rec pow x n =  
  if n < 1 then 1  
  else x * pow x (n - 1)
```

The switch to OCaml involves several significant issues:

1. The type of *pow* changes to $int \rightarrow int \rightarrow int$ since OCaml has no special type for \mathbb{N} (as is typical for execution-oriented programming languages). Thus the OCaml function admits arguments that are not admissible for the mathematical function. We speak of **spurious arguments**.
2. To make *pow* terminating for negative n , we return 1 for all $n < 1$. We can also use the equivalent comparison $n \leq 0$. If we don't scare away from nontermination for spurious arguments, we can also use the equality test $n = 0$.
3. The OCaml type *int* doesn't give us the full mathematical type \mathbb{Z} of integers but just a finite interval of machine integers.

When we design a function, there is always a mathematical level governing the coding level for OCaml. One important point about the mathematical level is that it doesn't change when we switch to another programming language. In this text, we will concentrate on the mathematical level.

When we argue the *correctness* of the function *pow*, we do this at the mathematical level using the infinite types \mathbb{Z} and \mathbb{N} . As it comes to the realization in OCaml, we just hope that the numbers involved for particular examples are small enough so that the difference between mathematical arithmetic and *machine arithmetic* doesn't show. The reason we ignore machine integers at the mathematical level is simplicity.

1.7 More about Mathematical Functions

We use the opportunity and give mathematical definitions for some functions we already discussed. A mathematical function definition consists of a type and a system of defining equations.

Remainder

$$\begin{aligned} \% : \mathbb{N} &\rightarrow \mathbb{N}^+ \rightarrow \mathbb{N} \\ x \% y &:= x && \text{if } x < y \\ x \% y &:= (x - y) \% y && \text{if } x \geq y \end{aligned}$$

Recall that \mathbb{N}^+ is the type of positive integers $1, 2, 3, \dots$. By using the type \mathbb{N}^+ for the divisor we avoid a division by zero.

1 Getting Started

Digit sum

$$\begin{aligned} D &: \mathbb{N} \rightarrow \mathbb{N} \\ D(x) &:= x && \text{if } x < 10 \\ D(x) &:= D(x/10) + (x \% 10) && \text{if } x \geq 10 \end{aligned}$$

Digit Reversal

$$\begin{aligned} R &: \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N} \\ R \ 0 \ a &:= a \\ R \ x \ a &:= R(x/10) (10 \cdot a + (x \% 10)) && \text{if } x > 0 \end{aligned}$$

A system of defining equations for a function must respect the type specified for the function. In particular, the defining equations must be *exhaustive* and *disjoint* for the type specified for the function.

Often, the defining equations of a function will be *terminating* for all arguments. This is the case for each of the three function defined above. In each case, the same termination argument applies: The first argument, which is a natural number, is decreased by each recursion step.

Curly braces

We can use curly braces to write several defining equations with the same left-hand side with just one left-hand side. For instance, we may write the defining equations of the digit sum function as follows:

$$D \ x \ := \ \begin{cases} x & \text{if } x < 10 \\ D(x/10) + (x \% 10) & \text{if } x \geq 10 \end{cases}$$

Total and partial functions

Functions where the defining equations terminate for all arguments are called **total**, and function where this is not the case are called **partial**. Most mathematical functions we will look at are total, but there are a few partial functions that matter. If there is a way to revise the mathematical definition of a function such that the function becomes total, we will usually do so. The reason we prefer total functions over partial functions is that equational reasoning for total functions is much easier than it is for partial functions. For correct equational reasoning about partial functions we need side conditions making sure that all function applications occurring in an equation used for reasoning do terminate.

1 Getting Started

We say that a function **diverges** for an argument if it does not terminate for this argument. Here is a function that diverges for all numbers smaller than 10 and terminates for all other numbers:

$$\begin{aligned} f &: \mathbb{N} \rightarrow \mathbb{N} \\ f(n) &:= f(n) && \text{if } n < 10 \\ f(n) &:= n && \text{if } n \geq 10 \end{aligned}$$

Exercise 1.7.1 Define a function $\mathbb{N} \rightarrow \mathbb{N}$ that terminates for even numbers and diverges for odd numbers.

Graph of a function

Abstractly, we may see a function $f : X \rightarrow Y$ as the set of all argument-result pairs $(x, f(x))$ where x is taken from the argument type X . We speak of the **graph** of a function. The graph of a function completely forgets about the definition of the function.

When we speak of a mathematical function in this text, we always include the definition of the function with a type and a system of defining equations. This information is needed so that we can compute with the function.

In noncomputational mathematics, one usually means by a function just a set of argument-result pairs.

GCD function

It is interesting to look at the mathematical version of the gcd function from §1.5:

$$\begin{aligned} G &: \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N} \\ G \ x \ 0 &:= x \\ G \ x \ y &:= G \ y \ (x \% y) && \text{if } y > 0 \end{aligned}$$

The defining equations are terminating since the second argument is decreased upon recursion (since $x \% y < y$ if $y > 0$). Note that the type of G admits $x = y = 0$, a case where no greatest common divisor exists.

1.8 A Higher-Order Function for Linear Search

A *boolean test for numbers* is a function $f : \text{int} \rightarrow \text{bool}$ expressing a condition for numbers. If $f(k) = \text{true}$, we say that k *satisfies* f .

Linear search is an algorithm that given a boolean test $f : \text{int} \rightarrow \text{bool}$ and a number n computes the first number $k \geq n$ satisfying f by checking f for

$$k = n, n + 1, n + 2, \dots$$

1 Getting Started

until f is satisfied for the first time. We realize linear search with an OCaml function

$$first : (int \rightarrow bool) \rightarrow int \rightarrow int$$

taking the test as first argument:

```
let rec first f k =  
  if f k then k  
  else first f (k + 1)
```

Functions taking functions as arguments are called **higher-order functions**. Higher-order functions are a key feature of functional programming.

Recall that we have characterized in §1.5 the integer quotient x/y as the maximal number k such that $k \cdot y \leq x$. Equivalently, we may characterize x/y as the first number $k \geq 0$ such that $(k + 1) \cdot y > x$ (recall the shelf interpretation). This gives us the equation

$$x/y = first (\lambda k. (k + 1) \cdot y > x) 0 \quad \text{if } x \geq 0 \text{ and } y > 0 \quad (1.1)$$

The functional argument of *first* is described with a **lambda expression**

$$\lambda k. (k + 1) \cdot y > x$$

Lambda expressions are a common mathematical notation for describing functions without giving them a name.⁶

We use Equation 1.1 to declare a function

$$div : int \rightarrow int \rightarrow int$$

computing quotients x/y :

```
let div x y = first (fun k -> (k + 1) * y > x) 0
```

From the declaration we learn that OCaml writes lambda expressions with the words “**fun**” and “**->**”. Here is a trace showing how quotients x/y are computed with *first*:

$$\begin{aligned} div\ 11\ 3 &= first\ (\lambda k. (k + 1) \cdot 3 > 11)\ 0 && 1 \cdot 3 \leq 11 \\ &= first\ (\lambda k. (k + 1) \cdot 3 > 11)\ 1 && 2 \cdot 3 \leq 11 \\ &= first\ (\lambda k. (k + 1) \cdot 3 > 11)\ 2 && 3 \cdot 3 \leq 11 \\ &= first\ (\lambda k. (k + 1) \cdot 3 > 11)\ 3 && 4 \cdot 3 > 11 \\ &= 3 \end{aligned}$$

⁶The greek letter “ λ ” is pronounced “lambda”. Sometimes lambda expressions $\lambda x.e$ are written with the more suggestive notation $x \mapsto e$.

1 Getting Started

We remark that *first* is our first inherently partial function. For instance, the function

$$\text{first } (\lambda k. \text{false}) : \mathbb{N} \rightarrow \mathbb{N}$$

diverges for all arguments. More generally, the application *first f n* diverges whenever there is no $k \geq n$ satisfying *f*.

Exercise 1.8.1 Declare a function *div'* such that $\text{div } x \ y = \text{div}' \ x \ y \ 0$ by specializing *first* to the test $\lambda k. (k + 1) \cdot y > x$.

Exercise 1.8.2 Declare a function *sqr* : $\mathbb{N} \rightarrow \mathbb{N}$ such that $\text{sqr}(n^2) = n$ for all *n*. Hint: Use *first*.

Exercise 1.8.3 Declare a terminating function *bounded_first* such that *bounded_first f n* yields the first $k \geq 0$ such that $k \leq n$ and *k* satisfies *f*.

1.9 Partial Applications

Functions described with lambda expressions can also be expressed with declared functions. To have an example, we declare *div* with *first* and a helper function *test* replacing the lambda expression:

```
let test x y k = (k + 1) * y > x
let div x y = first (test x y) 0
```

The type of the helper function *test* is $\text{int} \rightarrow \text{int} \rightarrow \text{int} \rightarrow \text{bool}$. Applying *test* to *x* and *y* yields a function of type $\text{int} \rightarrow \text{bool}$ as required by *first*. We speak of a **partial application**. Here is a trace showing how quotients x/y are computed with *first* and *test*:

$$\begin{aligned} \text{div } 11 \ 3 &= \text{first } (\text{test } 11 \ 3) \ 0 & 1 \cdot 3 &\leq 11 \\ &= \text{first } (\text{test } 11 \ 3) \ 1 & 2 \cdot 3 &\leq 11 \\ &= \text{first } (\text{test } 11 \ 3) \ 2 & 3 \cdot 3 &\leq 11 \\ &= \text{first } (\text{test } 11 \ 3) \ 3 & 4 \cdot 3 &> 11 \\ &= 3 \end{aligned}$$

We may describe the partial applications of *test* with equivalent lambda expressions:

$$\begin{aligned} \text{test } x \ y &= \lambda k. (k + 1) \cdot y > x \\ \text{test } 11 \ 3 &= \lambda k. (k + 1) \cdot 3 > 11 \end{aligned}$$

1 Getting Started

We can also describe partial applications of *test* to a single argument with equivalent lambda expressions:

$$\begin{aligned} \text{test } x &= \lambda y k. (k + 1) \cdot y > x = \lambda y. \lambda k. (k + 1) \cdot y > x \\ \text{test } 11 &= \lambda y k. (k + 1) \cdot y > 11 = \lambda y. \lambda k. (k + 1) \cdot y > 11 \end{aligned}$$

Note that lambda expressions with two argument variables are notation for nested lambda expressions with single arguments. We can also describe the function *test* with a nested lambda expression:

$$\text{test} = \lambda x y k. (k + 1) \cdot y > x = \lambda x. \lambda y. \lambda k. (k + 1) \cdot y > x$$

Following the nesting of lambda expressions, we may see applications and function types with several arguments as nestings of applications and function types with single arguments:

$$\begin{aligned} e_1 \ e_2 \ e_3 &= (e_1 \ e_2) \ e_3 \\ t_1 \rightarrow t_2 \rightarrow t_3 &= t_1 \rightarrow (t_2 \rightarrow t_3) \end{aligned}$$

Note that applications group to the left and function types group to the right.

We have considered *equations between functions* in the discussion of partial applications of *test*. We consider two functions as equal if they agree on all arguments. So functions with very different definitions may be equal. In fact, two functions are equal if and only if they have the same graph. We remark that there is no algorithm deciding equality of functions in general.

Exercise 1.9.1

- a) Write $\lambda x y k. (k + 1) \cdot y > x$ as a nested lambda expression.
- b) Write *test* 11 3 10 as a nested application.
- c) Write $\text{int} \rightarrow \text{int} \rightarrow \text{int} \rightarrow \text{bool}$ as a nested function type.

Exercise 1.9.2 Express the one-argument functions described by the expressions x^2 , x^3 and $(x+1)^2$ with lambda expressions in mathematical notation. Translate the lambda expressions to expressions in OCaml and have them type checked. Do the same for the two-argument function described by the expression $x < k^2$.

Exercise 1.9.3 (Sum functions)

- a) Define a function $\mathbb{N} \rightarrow \mathbb{N}$ computing the sum $0 + 1 + 2 + \dots + n$ of the first n numbers.

1 Getting Started

- b) Define a function $\mathbb{N} \rightarrow \mathbb{N}$ computing the sum $0 + 1^2 + 2^2 + \dots + n^2$ of the first n square numbers.
- c) Define a function $sum : (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \mathbb{N} \rightarrow \mathbb{N}$ computing for a given function f the sum $f(0) + f(1) + f(2) + \dots + f(n)$.
- d) Give partial applications of the function sum from (c) providing specialized sum functions as asked for by (a) and (b).

1.10 Inversion of Strictly Increasing Functions

A function $f : \mathbb{N} \rightarrow \mathbb{N}$ is **strictly increasing** if

$$f(0) < f(1) < f(2) < \dots$$

Strictly increasing functions can be inverted using linear search. That is, given a strictly increasing function $f : \mathbb{N} \rightarrow \mathbb{N}$, we can construct a function $g : \mathbb{N} \rightarrow \mathbb{N}$ such that $g(f(n)) = n$ for all n . The construction is explained by the equation

$$first (\lambda k. f(k+1) > f(n)) 0 = n \quad (1.2)$$

which in turn gives us the equation

$$(\lambda x. first (\lambda k. f(k+1) > x) 0) (f(n)) = n \quad (1.3)$$

For a concrete example, let $f(n) := n^2$. Equation 1.3 tells us that

$$g(x) := first (\lambda k. (k+1)^2 > x) 0$$

is a function $\mathbb{N} \rightarrow \mathbb{N}$ such that $g(n^2) = n$ for all n . Thus we know that

$$sqrt\ x := first (\lambda k. (k+1)^2 > x) 0$$

computes integer square roots $\lfloor \sqrt{x} \rfloor$. For instance, we have $sqrt(1) = 1$, $sqrt(4) = 2$, and $sqrt(9) = 3$. The **floor operator** $\lfloor x \rfloor$ converts a real number x into the greatest integer $y \leq x$.

Exercise 1.10.1 Give a trace for $sqrt\ 10$.

Exercise 1.10.2 Declare a function $sqrt'$ such that $sqrt\ x = sqrt'\ x\ 0$ by specializing $first$ to the test $\lambda k. (k+1)^2 > x$.

Exercise 1.10.3 The **ceiling operator** $\lceil x \rceil$ converts a real number into the least integer y such that $x \leq y$.

- a) Declare a function computing rounded down cube roots $\lfloor \sqrt[3]{x} \rfloor$.
- b) Declare a function computing rounded up cube roots $\lceil \sqrt[3]{x} \rceil$.

1 Getting Started

Exercise 1.10.4 Let $y > 0$. Convince yourself that $\lambda x. x/y$ inverts the strictly increasing function $\lambda n. n \cdot y$.

Exercise 1.10.5 Declare inverse functions for the following functions:

- a) $\lambda n. n^3$
- b) $\lambda n. n^k$ for $k \geq 2$
- c) $\lambda n. k^n$ for $k \geq 2$

Exercise 1.10.6 Declare a function $inv : (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow (\mathbb{N} \rightarrow \mathbb{N})$ that given a strictly increasing function f yields a function inverting f . Then express the functions from Exercise 1.10.5 using inv .

Exercise 1.10.7 Let $f : \mathbb{N} \rightarrow \mathbb{N}$ be strictly increasing. Convince yourself that the functions

$$\begin{aligned}\lambda x. first(\lambda k. f(k) = x) 0 \\ \lambda x. first(\lambda k. f(k) \geq x) 0 \\ \lambda x. first(\lambda k. f(k+1) > x) 0\end{aligned}$$

all invert f and find out how they differ.

1.11 Tail Recursion

A special form of functional recursion is **tail recursion**. Tail recursion matters since it can be executed more efficiently than general recursion. Tail recursion imposes the restriction that recursive function applications can only appear in tail positions where they directly yield the result of the function. Hence recursive applications appearing as part of another application (operator or function) are not tail recursive. Typical examples of tail recursive functions are the functions rev' , gcd , my_mod , and $first$ we have seen before. Counterexamples for tail recursive functions are the recursive functions pow (the recursive application is nested into a product) and my_div (the recursive application is nested into a sum).

Tail recursive functions have the property that their execution can be traced in a simple way. For instance, we have the tail recursive trace

$$\begin{aligned}gcd\ 36\ 132 &= gcd\ 132\ 36 \\ &= gcd\ 36\ 24 \\ &= gcd\ 24\ 12 \\ &= gcd\ 12\ 0 \\ &= 12\end{aligned}$$

1 Getting Started

For functions where the recursion is not tail recursive, traces look more complicated, for instance

$$\begin{aligned} \text{pow } 2 \ 3 &= 2 \cdot \text{pow } 2 \ 2 \\ &= 2 \cdot (2 \cdot \text{pow } 2 \ 1) \\ &= 2 \cdot (2 \cdot (2 \cdot \text{pow } 2 \ 0)) \\ &= 2 \cdot (2 \cdot (2 \cdot 1)) \\ &= 8 \end{aligned}$$

In imperative programming languages tail recursive functions can be expressed with loops. While imperative languages are designed such that loops should be used whenever possible, functional programming languages are designed such that tail recursive functions are preferable over loops.

Often recursive functions that are not tail recursive can be reformulated as tail recursive functions by introducing an extra argument serving as accumulator argument. Here is a tail recursive version of *pow*:

```
let rec pow' x n a =  
  if n < 1 then a  
  else pow' x (n - 1) (x * a)
```

We explain the role of the accumulator argument with a trace:

$$\begin{aligned} \text{pow}' \ 2 \ 3 \ 1 &= \text{pow}' \ 2 \ 2 \ 2 \\ &= \text{pow}' \ 2 \ 1 \ 4 \\ &= \text{pow}' \ 2 \ 0 \ 8 \\ &= 8 \end{aligned}$$

Exercise 1.11.1 (Factorials) In mathematics, the **factorial** of a positive integer n , denoted by $n!$, is the product of all positive integers less than or equal to n :

$$n! = n \cdot (n - 1) \cdot \dots \cdot 2 \cdot 1$$

For instance,

$$5! = 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1 = 120$$

In addition, $0!$ is defined as 1. We capture this specification with a recursive function defined as follows:

$$\begin{aligned} ! : \mathbb{N} &\rightarrow \mathbb{N} \\ 0! &:= 1 \\ (n + 1)! &:= (n + 1) \cdot n! \end{aligned}$$

1 Getting Started

- a) Declare a function $fac : int \rightarrow int$ computing factorials.
- b) Define a tail recursion function $f : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$ such that $n! = f\ 1\ n$.
- c) Declare a tail recursive function $fac' : int \rightarrow int \rightarrow int$ such that $fac'\ 1$ computes factorials.

1.12 Tuples

Sometimes we want a function that returns more than one value. For instance, the time for a marathon may be given with three numbers in the format $h : m : s$, where h is the number of hours, m is the number of minutes, and s is the number of seconds the runner needed. The time of Eliud Kipchoge's world record in 2018 in Berlin was $2 : 01 : 39$. There is the constraint that $m < 60$ and $s < 60$.

OCaml has **tuples** to represent collections of values as single values. To represent the marathon time of Kipchoge in Berlin, we can use the tuple

$$(2, 1, 39) : int \times int \times int$$

consisting of three integers. The product symbol “ \times ” in the **tuple type** is written as “ $*$ ” in OCaml. The **component types** of a tuple are not restricted to int , and there can be $n \geq 2$ **positions**, where n is called the **length** of the tuple. We may have tuples as follows:

$$\begin{aligned}(2, 3) &: int \times int \\ (7, true) &: int \times bool \\ ((2, 3), (7, true)) &: (int \times int) \times (int \times bool)\end{aligned}$$

Note that the last example nests tuples into tuples. We mention that tuples of length 2 are called **pairs**, and that tuples of length 3 are called **triples**.

We can now write two functions

$$\begin{aligned}sec &: int \times int \times int \rightarrow int \\ hms &: int \rightarrow int \times int \times int\end{aligned}$$

translating between times given in total seconds and times given as (h, m, s) tuples:

```
let sec (h,m,s) = 3600 * h + 60 * m + s
let hms x =
  let h = x / 3600 in
  let m = (x mod 3600) / 60 in
  let s = x mod 60 in
  (h,m,s)
```

1 Getting Started

Exercise 1.12.1

- a) Give a tuple of length 5 where the components are the values 2 and 3.
- b) Give a tuple of type $\text{int} \times (\text{int} \times (\text{bool} \times \text{bool}))$.
- c) Give a pair whose first component is a pair and whose second component is a triple.

Exercise 1.12.2 (Sorting triples) Declare a function *sort* sorting triples. For instance, we want $\text{sort } (3, 2, 1) = (1, 2, 3)$. Designing such a function is interesting. Given a triple (x, y, z) , the best solution we know of first ensures $y \leq z$ and then inserts x at the correct position. Start from the code snippet

```
let sort (x,y,z) =  
  let (y,z) = if y <= z then (y,z) else (z, y) in  
  if x <= y then ...?  
  else ...?
```

where the local declaration ensures $y \leq z$ using shadowing.

Exercise 1.12.3 (Medians) The median of three numbers is the number in the middle. For instance, the median of 5, 0, 1 is 1. Declare a function that takes three numbers and yields the median of the numbers.

1.13 Exceptions and Spurious Arguments

What happens when we execute the native operation $5/0$ in OCaml? Execution is aborted and an **exception** is reported:

Exception: Division_by_zero

Exceptions can be useful when debugging erroneous programs. We will say more about strings and exceptions in later chapters.

There is no equivalent to exceptions at the mathematical level. At the mathematical level we use types like \mathbb{N}^+ or side conditions like $y \neq 0$ to exclude undefined applications like $x/0$.

When coding a mathematical function in OCaml, we need to replace mathematical types like \mathbb{N} with the OCaml type *int*. This introduces **spurious arguments** not anticipated by the mathematical function. There are different ways to cope with spurious arguments:

1. Ignore the presence of spurious arguments. This is the best strategy when you solve exercises in this text.
2. Use a wrapper function raising exceptions when spurious arguments show up. The wrapper function facilitates the discovery of situations where functions are accidentally applied to spurious arguments.

1 Getting Started

As an example, we consider the coding of the mathematical remainder function

$$\begin{aligned} \text{rem} : \mathbb{N} \rightarrow \mathbb{N}^+ \rightarrow \mathbb{N} \\ \text{rem } x \ y &:= x && \text{if } x < y \\ \text{rem } x \ y &:= \text{rem } (x - y) \ y && \text{if } x \geq y \end{aligned}$$

as the OCaml function

```
let rec rem x y = if x < y then x else rem (x - y) y
```

receiving the type $\text{int} \rightarrow \text{int} \rightarrow \text{int}$. In OCaml we now have the spurious situation that $\text{rem } x \ 0$ diverges for all $x \geq 0$. There other spurious situations whose analysis is tedious since machine arithmetic needs to be taken into account. Using the wrapper function

```
let rem_checked x y =  
  if x >= 0 && y > 0 then rem x y  
  else invalid_arg "rem_checked"
```

all spurious situations are uniformly handled by throwing the exception

```
Invalid_argument "rem_checked"
```

There are several new features here:

- The lazy boolean and connective $x \geq 0 \ \&\& \ y > 0$ tests two conditions and is equivalent to `if x >= 0 then y > 0 else false`.
- There is the string `"rem_checked"`. Strings are values like integers and booleans and have type *string*.
- The predefined function `invalid_arg` raises an exception saying that `rem_checked` was called with spurious arguments.⁷

When an exception is raised, execution of a program is aborted and the exception raised is reported.

We use the opportunity and introduce the **lazy boolean connectives** as abbreviations for conditionals:

$$\begin{array}{llll} e_1 \ \&\& \ e_2 & \rightsquigarrow & \text{IF } e_1 \text{ THEN } e_2 \text{ ELSE false} & \text{lazy and} \\ e_1 \ || \ e_2 & \rightsquigarrow & \text{IF } e_1 \text{ THEN true ELSE } e_2 & \text{lazy or} \end{array}$$

Exercise 1.13.1 Consider the declaration

```
let eager_or x y = x || y
```

Find expressions e_1 and e_2 such that the expressions $e_1 || e_2$ and `eager_or e1 e2` behave differently. Hint: Choose a diverging expression for e_2 and keep in mind that execution of a function application

⁷OCaml says “invalid argument” for “spurious argument”.

1 Getting Started

executes all argument expressions. In contrast, execution of a conditional IF e_1 THEN e_2 ELSE e_3 executes e_1 and then either e_2 or e_3 , but not both.

Exercise 1.13.2 (Sorting triples) Recall Exercise 1.12.2. With lazy boolean connectives a function sorting triples can be written without much thinking by doing a naive case analysis considering the alternatives x is in the middle or y is in the middle or z is in the middle.

1.14 Summary

After working through this chapter you should be able to design and code functions computing powers, integer quotients and remainders, digit sums, digit reversals, and integer roots. You should understand that the design of functions happens at a mathematical level using mathematical types and equations. A given design can then be refined into a program in a given programming language. In this text we are using OCaml as programming language, assuming that this is the first programming language you see.⁸

You also saw a first higher-order function *first*, which can be used to obtain integer quotients and integer roots with a general scheme known as linear search. You will see many more examples of higher-order functions expressing basic computational schemes.

The model of computation we have assumed in this chapter is rewriting with defining equations. In this model, recursion appears in a natural way. You will have noticed that recursion is the feature where things get really interesting. We also have discussed tail recursion, a restricted form of recursion with nice properties we will study more carefully as we go on. All recursive functions we have seen in this chapter have equivalent tail recursive formulations (often using an accumulator argument).

Finally, we have seen tuples, which are compound values combining several values into a single value.

⁸Most readers will have done some programming in some programming language before starting with this text. Readers of this group often face the difficulty that they invest too much energy on mapping back the new things they see here to the form of programming they already understand. Since functional programming is rather different from other forms of programming, it is essential that you open yourself to the new ideas presented here. Keep in mind that a good programmer quickly adapts to new ways of thinking and to new programming languages.

2 Polymorphic Functions and Iteration

So far, some functions don't have unique types, as for instance the projection functions for pairs. The problem is solved by giving such functions schematic types called polymorphic types.

Using polymorphic types, we can declare a tail-recursive higher-order function iterating a function on a given value. It turns out that many functions can be obtained as instances of iteration. Our key example will be a function computing the n th prime number.

2.1 Polymorphic Functions

Consider the declaration of a projection function for pairs:

```
let fst (x,y) = x
```

What type does *fst* have? Clearly, $int \times int \rightarrow int$ and $bool \times int \rightarrow bool$ are both types admissible for *fst*. In fact, every type $t_1 \times t_2 \rightarrow t_1$ is admissible for *fst*. Thus there are infinitely many types admissible for *fst*.

OCaml solves the situation by typing *fst* with the polymorphic type

$$\forall \alpha \beta. \alpha \times \beta \rightarrow \alpha$$

A **polymorphic type** is a type scheme whose quantified variables (α and β in the example) can be instantiated with all types. The **instances** of the polymorphic type above are all types $t_1 \times t_2 \rightarrow t_1$ where the types t_1 and t_2 can be freely chosen. When a polymorphically typed identifier is used in an expression, it can be used with any instance of its polymorphic type. Thus *fst*(1, 2) and *fst*(true, 5) are both well-typed expressions.

Here is a polymorphic swap function for pairs:

```
let swap (x,y) = (y,x)
```

OCaml will type *swap* with the polymorphic type

$$\forall \alpha \beta. \alpha \times \beta \rightarrow \beta \times \alpha$$

This is in fact the **most general polymorphic type** that is admissible for *swap*. Similarly, the polymorphic type given for *fst* is the most general polymorphic type admissible for *fst*. OCaml will always derive most general types for function declarations.

2 Polymorphic Functions and Iteration

Exercise 2.1.1 Declare functions admitting the following polymorphic types:

- a) $\forall \alpha. \alpha \rightarrow \alpha$
- b) $\forall \alpha \beta. \alpha \rightarrow \beta \rightarrow \alpha$
- c) $\forall \alpha \beta \gamma. (\alpha \times \beta \rightarrow \gamma) \rightarrow \alpha \rightarrow \beta \rightarrow \gamma$
- d) $\forall \alpha \beta \gamma. (\alpha \rightarrow \beta \rightarrow \gamma) \rightarrow \alpha \times \beta \rightarrow \gamma$
- e) $\forall \alpha \beta. \alpha \rightarrow \beta$

2.2 Iteration

Given a function $f : t \rightarrow t$, we write $f^n(x)$ for the n -fold application of f to x . For instance, $f^0(x) = x$, $f^1(x) = f(x)$, $f^2(x) = f(f(x))$, and $f^3(x) = f(f(f(x)))$. More generally, we have

$$f^{n+1}(x) = f^n(fx)$$

Given an **iteration** $f^n(x)$, we call f the **step function** and x the **start value** of the iteration.

With iteration we can compute sums, products, and powers of non-negative integers just using additions $x + 1$:

$$\begin{aligned} x + n &= x + 1 \cdots + 1 = (\lambda a. a + 1)^n(x) \\ n \cdot x &= 0 + x \cdots + x = (\lambda a. a + x)^n(0) \\ x^n &= 1 \cdot x \cdots \cdot x = (\lambda a. a \cdot x)^n(1) \end{aligned}$$

Exploiting commutativity of addition and multiplication, we arrive at the defining equations

$$\begin{aligned} \text{succ } x &:= x + 1 \\ \text{add } x \ n &:= \text{succ}^n(x) \\ \text{mul } n \ x &:= (\text{add } x)^n(0) \\ \text{pow } x \ n &:= (\text{mul } x)^n(1) \end{aligned}$$

We define a polymorphic **iteration operator**

$$\begin{aligned} \text{iter} &: \forall \alpha. (\alpha \rightarrow \alpha) \rightarrow \mathbb{N} \rightarrow \alpha \rightarrow \alpha \\ \text{iter } f \ 0 \ x &:= x \\ \text{iter } f \ (n + 1) \ x &:= \text{iter } f \ n \ (f x) \end{aligned}$$

so that we can obtain iterations $f^n(x)$ as applications $\text{iter } f \ n \ x$ of the operator. Note that the function iter is polymorphic, higher-order, and tail-recursive. In OCaml, we will use the declaration

2 Polymorphic Functions and Iteration

```
let rec iter f n x =  
  if n < 1 then x  
  else iter f (n - 1) (f x)
```

Functions for addition, multiplication, and exponentiation can now be declared as follows:

```
let succ x = x + 1  
let add x y = iter succ y x  
let mul x y = iter (add y) x 0  
let pow x y = iter (mul x) y 1
```

Note that these declarations are non-recursive. Thus termination needs only be checked for *iter*, where it is obvious (2nd argument is decreased).

Exercise 2.2.1 Declare a function testing evenness of numbers by iterating on booleans. What do you have to change to obtain a function checking oddness?

Exercise 2.2.2 We have the equation

$$f^{n+1}(x) = f(f^n(x))$$

providing for an alternative, non-tail-recursive definition of an iteration operator. Give the mathematical definition and the declaration in OCaml of an iteration operator using the above equation.

2.3 Iteration on Pairs

Using iteration and successor as basic operations on numbers, we have defined functions computing sums, products, and powers of nonnegative numbers. We can also define a **predecessor function**¹

$$\begin{aligned} pred &: \mathbb{N}^+ \rightarrow \mathbb{N} \\ pred(n+1) &:= n \end{aligned}$$

just using iteration and successor (the successor of an integer x is $x+1$). The trick is to iterate on pairs. We start with the pair $(0,0)$ and iterate with a step function f such that $n+1$ iterations yield the pair $(n, n+1)$. For instance, the iteration

$$f^5(0,0) = f^4(0,1) = f^3(1,2) = f^2(2,3) = f(3,4) = (4,5)$$

using the **step function**

$$f(a,k) = (k, k+1)$$

¹the predecessor of an integer x is $x-1$.

2 Polymorphic Functions and Iteration

yields the predecessor of 5 as the first component of the computed pair. More generally, we have

$$(n, n + 1) = f^{n+1}(0, 0)$$

We can now declare a predecessor function as follows:

```
let pred n = fst (iter (fun (a,k) -> (k, succ k)) n (0,0))
```

Iteration on pairs is a powerful computation scheme. Our second example concerns the sequence of **Fibonacci numbers**

$$0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, \dots$$

which is well known in mathematics. The sequence is obtained by starting with 0, 1 and then adding new elements as the sum of the two preceding elements. We can formulate this method as a recursive function

$$\begin{aligned} fib &: \mathbb{N} \rightarrow \mathbb{N} \\ fib(0) &:= 0 \\ fib(1) &:= 1 \\ fib(n+2) &:= fib(n) + fib(n+1) \end{aligned}$$

Incidentally, this is the first recursive function we see where the recursion is **binary** (two recursive applications) rather than **linear** (one recursive application). Termination follows with the usual argument that each recursion step decreases the argument.

If we look again at the rule generating the Fibonacci sequence, we see that we can compute the sequence by starting with the pair (0, 1) and iterating with the step function $f(a, b) = (b, a + b)$. For instance,

$$f^5(0, 1) = f^4(1, 1) = f^3(1, 2) = f^2(2, 3) = f(3, 5) = (5, 8)$$

yields the pair $(fib(5), fib(6))$. More generally, we have

$$(fib(n), fib(n+1)) = (\lambda(a, b). (b, a + b))^n (0, 1)$$

Thus we can declare an iterative Fibonacci function as follows:

```
let fibi n = fst (iter (fun (a,b) -> (b, a + b)) n (0,1))
```

In contrast to the previously defined function *fib*, function *fibi* requires only tail recursion as provided by *iter*.

Exercise 2.3.1 Declare a function computing the sum $0+1+2+\dots+n$ by iteration starting from the pair (0, 1).

2 Polymorphic Functions and Iteration

Exercise 2.3.2 Declare a function $f : \mathbb{N} \rightarrow \mathbb{N}$ computing the sequence

$$0, 1, 1, 2, 4, 7, 13, \dots$$

obtained by starting with 0, 1, 1 and then adding new elements as the sum of the three preceding elements. For instance, $f(3) = 2$, $f(4) = 4$, and $f(5) = 7$.

Exercise 2.3.3 Functions defined with iteration can always be elaborated into tail-recursive functions not using iteration. If the iteration is on pairs, one can use separate accumulator arguments for the components of the pairs. Follow this recipe and declare a tail-recursive function fib' such that $fib' \ n \ 0 \ 1 = fib(n)$.

Exercise 2.3.4 Recall the definition of factorials $n!$ from Exercise 1.11.1.

- a) Give a step function f such that $(n!, n) = f^n(1, 0)$.
- b) Declare a function $faci$ computing factorials with iteration.
- c) Declare a tail-recursive function fac' such that $fac' \ n \ 1 \ 0 = n!$. Follow the recipe from Exercise 2.3.3.

2.4 Computing Primes

A **prime number** is an integer greater 1 that cannot be obtained as the product of two integers greater 1. There are infinitely many prime numbers. The sequence of prime numbers starts with

$$2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41, 43, 47, \dots$$

We want to declare a function $nth_prime : \mathbb{N} \rightarrow \mathbb{N}$ that yields the elements of the sequence starting with $nth_prime \ 0 = 2$. Assuming a primality test $prime : int \rightarrow bool$, we can declare nth_prime as follows:

```
let next_prime x = first prime (x + 1)
let nth_prime n = iter next_prime n 2
```

Note that $next_prime \ x$ yields the first prime greater than x . Also note that the elegant declarations of $next_prime$ and nth_prime are made possible by the higher-order functions *first* and *iter*.

It remains to come up with a **primality test** (a test checking whether an integer is a prime number). Here are 4 equivalent characterizations of primality of x assuming that x , k , and n are natural numbers:

2 Polymorphic Functions and Iteration

1. $x \geq 2 \wedge \forall k \geq 2. \forall n \geq 2. x \neq k \cdot n$
2. $x \geq 2 \wedge \forall k, 2 \leq k < x. x \% k > 0$
3. $x \geq 2 \wedge \forall k > 1, k^2 \leq x. x \% k > 0$
4. $x \geq 2 \wedge \forall k, 1 < k \leq \sqrt[3]{x}. x \% k > 0$

Characterization (1) is close to the informal definition of prime numbers. Characterizations (2), (3), and (4) are useful for our purposes since they can be realized algorithmically. Starting from (1), we see that $x - 1$ can be used as an upper bound for k and n . Thus we have to test $x = k \cdot n$ only for finitely many k and n , which can be done algorithmically. The next thing we see is that it suffices to have k because n can be kept implicit using the remainder operation. Finally, we see that it suffices to test for k such that $k^2 \leq x$ since we can assume $n \geq k$. Thus we can sharpen the upper bound from $x - 1$ to \sqrt{x} .

Here we choose the second characterization to declare a primality test in OCaml and leave the realization of the computationally faster fourth characterization as an exercise.

To test the bounded universal quantification in the second characterization, we declare a higher-order function

$$forall : int \rightarrow int \rightarrow (int \rightarrow bool) \rightarrow bool$$

such that

$$forall\ m\ n\ f = \text{true} \iff \forall k, m \leq k \leq n. f\ k = \text{true}$$

using the defining equations²

$$forall\ m\ n\ f := \begin{cases} \text{true} & m > n \\ f\ m \ \&\& \ forall\ (m+1)\ n\ f & m \leq n \end{cases}$$

The function *forall* terminates since $|n + 1 - m| \geq 0$ decreases with every recursion step.

We now define a primality test based on the second characterization:

$$prime\ x := x \geq 2 \ \&\& \ forall\ 2\ (x - 1)\ (\lambda k. x \% k > 0)$$

The discussion of primality tests makes it very clear that programming involves mathematical reasoning.

Efficient primality tests are important in cryptography and other areas of computer science. The naive versions we have discussed here are known as trial division algorithms and are too slow for practical purposes.

²We use the curly brace as an abbreviation to write two equations as a single equation.

2 Polymorphic Functions and Iteration

Exercise 2.4.1 Declare a test $exists : int \rightarrow int \rightarrow (int \rightarrow bool) \rightarrow bool$ such that $exists\ m\ n\ f = \text{true} \iff \exists k, m \leq k \leq n. f\ k = \text{true}$ in two ways:

- a) Directly following the design of *forall*.
- b) Using *forall* and boolean negation $not : bool \rightarrow bool$.

Exercise 2.4.2 Declare a primality test based on the fourth characterization (i.e., upper bound $\sqrt[3]{x}$). Convince yourself with an OCaml interpreter that testing with upper bound $\sqrt[3]{x}$ is much faster on large primes (check 479,001,599 and 87,178,291,199).

Exercise 2.4.3 Explain why the following functions are primality tests:

- a) $\lambda x. x \geq 2 \ \&\& \ first\ (\lambda k. x \% k = 0)\ 2 = x$
- b) $\lambda x. x \geq 2 \ \&\& \ (first\ (\lambda k. k^2 \geq x \ || \ x \% k = 0)\ 2)^2 > x$

Hint for (b): Let k be the number the application of *first* yields. Distinguish three cases: $k^2 < x$, $k^2 = x$, and $k^2 > x$.

Exercise 2.4.4 Convince yourself that the four characterizations of primality given above are equivalent.

2.5 Polymorphic Typing Rules

We distinguish between **monomorphic** and **polymorphic types**. Polymorphic types (e.g., $\forall \alpha. \alpha \rightarrow \alpha$) are type schemes where one or several type variables are universally quantified. Monomorphic types (e.g., $t \rightarrow t$) may contain type variables but must not contain quantifiers. Polymorphic types are obtained from monomorphic types by quantifying one or several variables. This disallows nesting of polymorphic types into function types or tuple types.

The typing rules from §?? can be extended to polymorphic types. We will write t for monomorphic types and T for polymorphic types. A type environment E can now contain both monomorphic bindings $x : t$ and polymorphic bindings $x : T$. The rules from §?? stay unchanged except for the now more general E and t . Three new **polymorphic typing rules** are added to handle polymorphic types.

The first polymorphic typing rule provides for the use of polymorphic variables in expressions:

$$\frac{(x : T) \in E \quad T \succ t}{E \vdash x : t}$$

2 Polymorphic Functions and Iteration

The notation $T \succ t$ says that the monomorphic type t is an instance of the polymorphic type T . We can now derive the typing

$$f : \forall \alpha. \alpha \rightarrow \alpha \vdash f : \text{int} \rightarrow \text{int}$$

The second polymorphic typing rule accommodates let expressions declaring polymorphic functions:

$$\frac{E \vdash e_1 : t_1 \quad E, x : T \vdash e_2 : t \quad T \dot{\prec}^E t_1 \quad e_1 = \lambda \dots}{E \vdash \text{LET } x = e_1 \text{ IN } e_2 : t}$$

The side condition $T \dot{\prec}^E t$ says that the monomorphic type t is a most general instance of the polymorphic type T obtained by erasing the quantifier prefix of T after the quantified type variables have been renamed so that they don't occur in E . The side condition $e_1 = \lambda \dots$ says that e_1 must be a lambda expression. We can now derive the typing

$$\square \vdash \text{LET } f = \lambda x. x \text{ IN } (f\ 5, f\ \text{true}) : \text{int} \times \text{bool}$$

where f is assigned the polymorphic type $\forall \alpha. \alpha \rightarrow \alpha$ and the most general instance $\alpha \rightarrow \alpha$ by the rule for let expressions. The two using occurrences of f are typed with the instances $\text{int} \rightarrow \text{int}$ and $\text{bool} \rightarrow \text{bool}$ of $\forall \alpha. \alpha \rightarrow \alpha$.

Note that the two rules for polymorphic types are not algorithmic since they involve guessing of types. It turns out that OCaml comes with a smart type inference algorithm doing all type guessing in an algorithmic way. In fact you can use OCaml to find out for given E and e the most general type t such that $E \vdash e : t$.

Example: Polymorphic self application

As an example we consider the expression

$$\text{LET } f = \lambda x. x \text{ IN } ff$$

which defines f as an identity function. This expression can only be typed if f is given a polymorphic type since $f : t \rightarrow t$ cannot be applied to itself (since $t \neq t \rightarrow t$). However, the polymorphic typing rule for let expressions makes it possible to assign to f the polymorphic type $\forall \alpha. \alpha \rightarrow \alpha$, which makes it possible to type the self application ff with the instances $(\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$ and $\alpha \rightarrow \alpha$. Here is a derivation establishing a most general type for ff (the derivation is given in two parts to fit on the page):

$$\frac{\frac{E, f : \forall \alpha. \alpha \rightarrow \alpha \vdash f : (\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)}{E, f : \forall \alpha. \alpha \rightarrow \alpha \vdash f : \alpha \rightarrow \alpha} \quad \frac{E, f : \forall \alpha. \alpha \rightarrow \alpha \vdash f : \alpha \rightarrow \alpha}{E, f : \forall \alpha. \alpha \rightarrow \alpha \vdash ff : \alpha \rightarrow \alpha}}$$

2 Polymorphic Functions and Iteration

$$\frac{\overline{E, x : \alpha \vdash x : \alpha}}{E \vdash \lambda x. x : \alpha \rightarrow \alpha} \quad \frac{\dots \quad \dots}{E, f : \forall \alpha. \alpha \rightarrow \alpha \vdash ff : \alpha \rightarrow \alpha} \\ \hline E \vdash \text{LET } f = \lambda x. x \text{ IN } ff : \alpha \rightarrow \alpha$$

The derivation assumes that E does not contain the type variable α .

Exercise 2.5.1 (Untypable expressions) We say that an expression e is **typable** if there are an environment E and a type t such that the judgment $E \vdash e : t$ is derivable. Explain why the following expression are not typable.

a) $\lambda f.f f$

b) $(\lambda f. (f\ 5, f\ "a")) (\lambda x.x)$

Polymorphic typing rule for recursive let expressions

We also need a polymorphic typing rule for recursive let expressions. This rule doesn't need new ideas but can be obtained as a combination of the monomorphic typing rule for recursive let expressions and the polymorphic typing rule for let expressions we have just seen:

$$\frac{E, f : t_1 \rightarrow t_2, x : t_1 \vdash e_1 : t_2 \quad E, f : T \vdash e_2 : t \quad T \succ^E t_1 \rightarrow t_2}{E \vdash \text{LET REC } f x = e_1 \text{ IN } e_2 : t}$$

Note that f is typed monomorphically for e_1 and polymorphically for e_2 .³

Keep in mind that only variables declared by `let` expressions can be assigned polymorphic types.

2.6 Polymorphic Exception Raising and Equality Testing

Recall the predefined function *invalid_arg* discussed in §1.13. Like every function in OCaml, *invalid_arg* must be accommodated with a type. It turns out that the natural type for *invalid_arg* is a polymorphic type:

$$invalid \quad arg : \forall \alpha. string \rightarrow \alpha$$

With this type an application of *invalid_arg* can be typed with whatever type is required by the context of the application. Since evaluation of the application raises an exception and doesn't yield a value, the return type of the function doesn't matter for evaluation.

³In principle, f can be typed polymorphically for both constituents e_1 and e_2 of a recursive let expression, but this generalization of the polymorphic typing rule for recursive let expressions ruins the type inference algorithm for polymorphic types.

2 Polymorphic Functions and Iteration

Like functions, operations must be accommodated with types. So what is the type of the equality test? OCaml follows common mathematical practice and admits the equality test for all types. To this purpose, the operator testing equality is accommodated with a polymorphic type:

$$= : \forall \alpha. \alpha \rightarrow \alpha \rightarrow \text{bool}$$

The polymorphic type promises more than the equality test delivers. There is the general problem that a meaningful equality test for functions cannot be realized computationally. OCaml bypasses the problem by the crude provision that an equality test on functions raises an invalid argument exception.

We assume that functions at the mathematical level always return values and do not raise exceptions. We handle division by zero by assuming that it returns some value, say 0. Similarly, we handle an equality test for functions by assuming that it always returns `true`. When reasoning at the mathematical level, we will avoid situations where the ad hoc definitions come into play, following common mathematical practice.

2.7 Summary

Functions can receive polymorphic types. Polymorphic types are type schemes whose variables serve as placeholders for arbitrary types. A polymorphic function can be used with every instance of its type scheme. Typical examples for polymorphic functions (i.e., functions that are assigned a polymorphic type) are functions for tuples and the iteration operator.

OCaml will assign polymorphic types to declared functions whenever this is possible. This way OCaml's type inference algorithm can infer unique polymorphic types for functions that don't have most general monomorphic types.

Type checking as explained in §?? extends to polymorphic types by adding three polymorphic typing rules. Only variables declared with `let` expressions can receive polymorphic types. In particular, lambda expressions cannot receive polymorphic types.

We remark that the dynamic semantics of programs is not affected by polymorphic types since it doesn't know about static types anyway. We may see polymorphic types as a means for making programs well typed that otherwise would not be well typed. An example is the iteration function which can iterate on different types (we have seen iterations on numbers, booleans, and pairs). Another example are graded projections for tuples, which can be used for different component types.

2 Polymorphic Functions and Iteration

Iteration is a tail-recursive computation scheme that can be formulated as a polymorphic higher-order function. Without functional arguments it is impossible to formulate iteration as a function, and without polymorphism we cannot formulate iteration for all iteration types.

3 Lists

Lists are a basic mathematical data structure providing a recursive representation for finite sequences. Lists are essential for programming. OCaml, and functional programming languages in general, accommodate lists in a mathematically clean way. Three interesting problems we will attack with lists are decimal representation, sorting, and prime factorization:

$$\begin{aligned} dec\ 735 &= [7, 3, 5] \\ sort\ [7, 2, 7, 6, 3, 4, 5, 3] &= [2, 3, 3, 4, 5, 6, 7, 7] \\ prime_fac\ 735 &= [3, 5, 7, 7] \end{aligned}$$

3.1 Nil and Cons

A list represents a finite sequence $[x_1, \dots, x_n]$ of values. All **elements** of a list must have the same type. A **list type** $\mathcal{L}(t)$ contains all lists whose elements are of type t ; we speak of **lists over** t . For instance,

$$\begin{aligned} [1, 2, 3] &: \mathcal{L}(\mathbb{Z}) \\ [\text{true}, \text{true}, \text{false}] &: \mathcal{L}(\mathbb{B}) \\ [(\text{true}, 1), (\text{true}, 2), (\text{false}, 3)] &: \mathcal{L}(\mathbb{B} \times \mathbb{Z}) \\ [[1, 2], [3], []] &: \mathcal{L}(\mathcal{L}(\mathbb{Z})) \end{aligned}$$

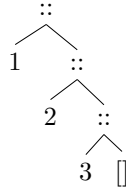
All lists are obtained from the **empty list** $[]$ using the binary constructor **cons** written as “ $::$ ”:

$$\begin{aligned} [1] &= 1 :: [] \\ [1, 2] &= 1 :: (2 :: []) \\ [1, 2, 3] &= 1 :: (2 :: (3 :: [])) \end{aligned}$$

The empty list $[]$ also counts as a constructor and is called **nil**. Given a nonempty list $x :: l$ (i.e., a list obtained with cons), we call x the **head** and l the **tail** of the list.

It is important to see lists as trees. For instance, the list $[1, 2, 3]$ may be depicted as the tree

3 Lists



The tree representation shows how lists are obtained with the constructors `nil` and `cons`. It is important to keep in mind that the bracket notation $[x_1, \dots, x_n]$ is just notation for a list obtained with n applications of `cons` from `nil`. Also keep in mind that every list is obtained with either the constructor `nil` or the constructor `cons`.

Notationally, `cons` acts as an infix operator grouping to the right. Thus we can omit the parentheses in $1 :: (2 :: (3 :: []))$. Moreover, we have $[x_1, \dots, x_n] = x_1 :: \dots :: x_n :: []$.

Given a list $[x_1, \dots, x_n]$, we call the values x_1, \dots, x_n the **elements** or the **members** of the list.

Despite the fact that tuples and lists both represent sequences, tuple types and list types are quite different:

- A tuple type $t_1 \times \dots \times t_n$ admits only tuples of length n , but may fix different types for different components.
- A list type $\mathcal{L}(t)$ admits lists $[x_1, \dots, x_n]$ of any length but fixes a single type t for all elements.

Exercise 3.1.1 Give the types of the following lists and tuples.

- | | | |
|----------------|-----------------------|-----------------------|
| a) $[1, 2, 3]$ | c) $[(1, 2), (2, 3)]$ | e) $[[1, 2], [2, 3]]$ |
| b) $(1, 2, 3)$ | d) $((1, 2), (2, 3))$ | |

3.2 Basic List Functions

The fact that all lists are obtained with `nil` and `cons` facilitates the definition of basic operations on lists. We start with the definition of a polymorphic function

$$\begin{aligned}
 \text{length} &: \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathbb{N} \\
 \text{length } [] &:= 0 \\
 \text{length } (x :: l) &:= 1 + \text{length } l
 \end{aligned}$$

3 Lists

that yields the **length** of a list. We have $length[x_1, \dots, x_n] = n$. Another prominent list operation is **concatenation**:

$$\begin{aligned} @ : \forall \alpha. \mathcal{L}(\alpha) &\rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \\ [] @ l_2 &:= l_2 \\ (x :: l_1) @ l_2 &:= x :: (l_1 @ l_2) \end{aligned}$$

We have $[x_1, \dots, x_m] @ [y_1, \dots, y_n] = [x_1, \dots, x_m, y_1, \dots, y_n]$. We call $l_1 @ l_2$ the **concatenation** of l_1 and l_2 .

How can we define an operation **reversing** lists:

$$rev[x_1, \dots, x_n] = [x_n, \dots, x_1]$$

For instance, $rev[1, 2, 3] = [3, 2, 1]$. To define rev , we need defining equations for nil and $cons$. The defining equation for nil is obvious, since the reversal of the empty list is the empty list. For $cons$ we use the equation $rev(x :: l) = rev(l) @ [x]$ which expresses a basic fact about reversal. This brings us to the following definition of **list reversal**:

$$\begin{aligned} rev : \forall \alpha. \mathcal{L}(\alpha) &\rightarrow \mathcal{L}(\alpha) \\ rev [] &:= [] \\ rev (x :: l) &:= rev(l) @ [x] \end{aligned}$$

We can also define a tail recursive list reversal function. As usual we need an accumulator argument. The resulting function combines reversal and concatenation:

$$\begin{aligned} rev_append : \forall \alpha. \mathcal{L}(\alpha) &\rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \\ rev_append [] l_2 &:= l_2 \\ rev_append (x :: l_1) l_2 &:= rev_append l_1 (x :: l_2) \end{aligned}$$

We have $rev(l) = rev_append l []$. The following trace shows the defining equations of rev_append at work:

$$\begin{aligned} rev_append [1, 2, 3] [] &= rev_append [2, 3] [1] \\ &= rev_append [3] [2, 1] \\ &= rev_append [] [3, 2, 1] \\ &= [3, 2, 1] \end{aligned}$$

The functions defined so far are all defined by **list recursion**.¹ List recursion means that there is no recursion for the empty list, and that

¹List recursion is better known as *structural recursion on lists*.

3 Lists

for every nonempty list $x :: l$ the recursion is on the tail l of the list. List recursion always terminates since lists are obtained from nil with finitely many applications of cons.

Another prominent list operation is

$$\begin{aligned} \text{map} &: \forall \alpha \beta. (\alpha \rightarrow \beta) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\beta) \\ \text{map } f \ [] &:= [] \\ \text{map } f (x :: l) &:= f x :: \text{map } f l \end{aligned}$$

We have $\text{map } f [x_1, \dots, x_n] = [f x_1, \dots, f x_n]$. We say that map applies the function f **pointwise** to a list. Note that map is defined once more with list recursion.

Finally, we define a function that yields the list of all integers between two numbers m and n (including m and n):

$$\begin{aligned} \text{seq} &: \mathbb{Z} \rightarrow \mathbb{Z} \rightarrow \mathcal{L}(\mathbb{Z}) \\ \text{seq } m \ n &:= \text{IF } m > n \text{ THEN } [] \text{ ELSE } m :: \text{seq } (m + 1) \ n \end{aligned}$$

For instance, $\text{seq } -1 \ 5 = [-1, 0, 1, 2, 4, 5]$. This time the recursion is not on a list but on the number m . The recursion terminates since every recursion step decreases $n - m$ and recursion stops once $n - m < 0$.

Exercise 3.2.1 Define a function $\text{null} : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathbb{B}$ testing whether a list is the empty list. Do not use the equality test.

Exercise 3.2.2 Consider the expression $1 :: 2 :: [] @ 3 :: 4 :: []$.

1. Put in the redundant parentheses.
2. Give the list the expression evaluates to in bracket notation.
3. Give the tree representation of the list the expression evaluates to.

Exercise 3.2.3 Decide for each of the following equations whether it is well-typed, and, in case it is well-typed, whether it is true. Assume $l_1, l_2 : \text{List}(\mathbb{N})$.

- a) $1 :: 2 :: 3 = 1 :: [2, 3]$
- b) $1 :: 2 :: 3 :: [] = 1 :: (2 :: [3])$
- c) $l_1 :: [2] = l_1 @ [2]$
- d) $(l_1 @ [2]) @ l_2 = l_1 @ (2 :: l_2)$
- e) $(l_1 :: 2) @ l_2 = l_1 @ (2 :: l_2)$
- f) $\text{map } (\lambda x. x^2) [1, 2, 3] = [1, 4, 9]$
- g) $\text{rev } (l_1 @ l_2) = \text{rev } l_2 @ \text{rev } l_1$

3.3 List Functions in OCaml

Given a mathematical definition of a list function, it is straightforward to declare the function in OCaml. The essential new construct are so-called **match expressions** making it possible to discriminate between empty and nonempty lists. Here is a declaration of a length function:

```
let rec length l =
  match l with
  | [] -> 0
  | x :: l -> 1 + length l
```

The match expression realizes a case analysis on lists using separate **rules** for the empty list and for nonempty lists. The left hand sides of the rules are called **patterns**. The cons pattern applies to nonempty lists and binds the **local variables** x and l to the head and the tail of the list.

Note that the pattern variable l introduced by the second rule of the match shadows the argument variable l in the declaration of *length*. The shadowing could be avoided by using a different variable name for the pattern variable (for instance, l').

Below are declarations of OCaml functions realizing the functions *append*, *rev_append*, and *map* defined before. Note how the defining equations translate into rules of match expressions.

```
let rec append l1 l2 =
  match l1 with
  | [] -> l2
  | x :: l1 -> x :: append l1 l2

let rec rev_append l1 l2 =
  match l1 with
  | [] -> l2
  | x :: l1 -> rev_append l1 (x :: l2)

let rec map f l =
  match l with
  | [] -> []
  | x :: l -> f x :: map f l
```

For each of the function declarations, OCaml infers the polymorphic type we have specified with the mathematical definition.

We remark that OCaml realizes the bracket notation for lists using semicolons to separate elements. For instance, the list $[1, 2, 3]$ is written as `[1; 2; 3]` in OCaml.

OCaml provides predefined functions for lists as **fields** of a predefined **standard module** *List*. For instance:

3 Lists

```

List.length      : 'a list -> int
List.append      : 'a list -> 'a list -> 'a list
List.rev_append  : 'a list -> 'a list -> 'a list
List.rev         : 'a list -> 'a list
List.map         : ('a -> 'b) -> 'a list -> 'b list

```

The above listing uses OCaml notation:

- **Dot notation** is used to name the fields of modules; for instance, *List.append* denotes the field *append* of the module *List*.
- List types $\mathcal{L}(t)$ are written in reverse order as “*t list*”.
- Type variables are written with a leading quote, for instance, “*a*”
- The quantification prefix of polymorphic types is suppressed, relying on the assumption that all occurring type variables are quantified.

Here are further notational details concerning lists in OCaml:

- List concatenation *List.append* is also available through the infix operator “@”.
- The infix operators “::” and “@” both group to the right, and “::” takes its arguments before “@”. For instance,

$$1 :: 2 :: [3; 4] @ [5] \rightsquigarrow (1 :: (2 :: [3; 4])) @ [5]$$

- The operators “::” and “@” take their arguments before comparisons and after arithmetic operations.

Exercise 3.3.1 Declare a function *seq* following the mathematical definition in the previous section.

Exercise 3.3.2 (Init) Declare a polymorphic function *init* such that $\text{init } n \ f = [f(0), \dots, f(n-1)]$ for $n \geq 0$. Note that n is the length of the result list. Write your function with a tail-recursive helper function. Make sure your function agrees with OCaml’s predefined function *List.init*. Use *List.init* to declare polymorphic functions that yield lists $[f(m), f(m+1), \dots, f(m+n-1)]$ and $[f(m), f(m+1), \dots, f(n)]$.

Exercise 3.3.3 Declare a function *flatten* : $\forall \alpha. \mathcal{L}(\mathcal{L}(\alpha)) \rightarrow \mathcal{L}(\alpha)$ concatenating the lists appearing as elements of a given list:

$$\text{flatten } [l_1, \dots, l_n] = l_1 @ \dots @ l_n @ []$$

For instance, we want $\text{flatten } [[1, 2], [], [3], [4, 5]] = [1, 2, 3, 4, 5]$.

Exercise 3.3.4 (Decimal Numbers) With lists we have a mathematical representation for decimal numbers. For instance, the decimal representation for the natural number 1234 is the list $[1, 2, 3, 4]$.

3 Lists

- a) Declare a function $dec : \mathbb{N} \rightarrow \mathcal{L}(\mathbb{N})$ that yields the decimal number for a natural number. For instance, we want $dec\ 1324 = [1, 3, 2, 4]$.
- b) Declare a function $num : \mathcal{L}(\mathbb{N}) \rightarrow \mathbb{N}$ that converts decimal numbers into numbers: $num(dec\ n) = n$.

Hint: Declare num with a tail-recursive function num' such that, for instance,

$$\begin{aligned} num\ [1, 2, 3] &= num'\ [1, 2, 3]\ 0 \\ &= num'\ [2, 3]\ 1 \\ &= num'\ [3]\ 12 \\ &= num'\ []\ 123 = 123 \end{aligned}$$

Exercise 3.3.5 Declare functions

$$\begin{aligned} zip &: \forall \alpha \beta. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\beta) \rightarrow \mathcal{L}(\alpha \times \beta) \\ unzip &: \forall \alpha \beta. \mathcal{L}(\alpha \times \beta) \rightarrow \mathcal{L}(\alpha) \times \mathcal{L}(\beta) \end{aligned}$$

such that

$$\begin{aligned} zip\ [x_1, \dots, x_n]\ [y_1, \dots, y_n] &= [(x_1, y_1), \dots, (x_n, y_n)] \\ unzip\ [(x_1, y_1), \dots, (x_n, y_n)] &= ([x_1, \dots, x_n], [y_1, \dots, y_n]) \end{aligned}$$

3.4 Fine Points About Lists

We speak of the collection of list types $\mathcal{L}(t)$ as a **type family**. We may describe the family of list types with the grammar

$$\mathcal{L}(\alpha) ::= [] \mid \alpha :: \mathcal{L}(\alpha)$$

fixing the **constructors** `nil` and `cons`. Speaking semantically, every value of a list type is obtained with either the constructor `nil` or the constructor `cons`. Moreover, `[]` is a value that is a member of every list type $\mathcal{L}(t)$, and $v_1 :: v_2$ is a value that is a member of a list type $\mathcal{L}(t)$ if the value v_1 is a member of the type t and the value v_2 is a member of the type $\mathcal{L}(t)$.

As it comes to type checking, the constructors `nil` and `cons` are accommodated with polymorphic types:

$$\begin{aligned} [] &: \forall \alpha. \mathcal{L}(\alpha) \\ (::) &: \forall \alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \end{aligned}$$

OCaml comes with the peculiarity that constructors taking arguments (e.g., `cons`) can only be used when applied to all arguments. We can

3 Lists

bypass this restriction by declaring a polymorphic function applying the `cons` constructor:

$$\begin{aligned} \text{cons} &: \forall \alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \\ \text{cons } x \ l &:= x :: l \end{aligned}$$

OCaml provides this function as `List.cons`.

3.5 Membership and List Quantification

We define a polymorphic function that tests whether a value appears as element of a list:

$$\begin{aligned} \text{mem} &: \forall \alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B} \\ \text{mem } _ \ [] &:= \text{false} \\ \text{mem } x \ (y :: l) &:= (x = y) \ || \ \text{mem } x \ l \end{aligned}$$

Note that `mem` tail-recurses on the list argument. Also recall the discussion of OCaml's polymorphic equality test in §2.6. We will write $x \in l$ to say that x is an element of the list l .

The structure of the membership test can be generalized with a polymorphic function that for a test and a list checks whether some element of the list satisfies the test:

$$\begin{aligned} \text{exists} &: \forall \alpha. (\alpha \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B} \\ \text{exists } p \ [] &:= \text{false} \\ \text{exists } p \ (x :: l) &:= p \ x \ || \ \text{exists } p \ l \end{aligned}$$

The expression

$$\text{exists } (\lambda x. x = 5) : \mathcal{L}(\mathbb{Z}) \rightarrow \mathbb{B}$$

now gives us a test that checks whether a lists of numbers contains the number 5. More generally, we have

$$\text{mem } x \ l = \text{exists } ((=)x) \ l$$

Exercise 3.5.1 Declare `mem` and `exists` in OCaml. For `mem` consider two possibilities, one with `exists` and one without helper function.

Exercise 3.5.2 Convince yourself that `exists` is tail-recursive (by eliminating the derived form `||`).

3 Lists

Exercise 3.5.3 Declare a tail-recursive function

$$forall : \forall \alpha. (\alpha \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B}$$

testing whether all elements of a list satisfy a given test. Consider two possibilities, one without a helper function, and one with *exists* exploiting the equivalence $(\forall x \in l. p(x)) \longleftrightarrow (\neg \exists x \in l. \neg p(x))$.

Exercise 3.5.4 Declare a function $count : \forall \alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{N}$ that counts how often a value appears in a list. For instance, we want $count\ 5\ [2, 5, 3, 5] = 2$.

Exercise 3.5.5 (Inclusion) Declare a function

$$incl : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B}$$

which tests whether all elements of the first list are elements of the second list.

Exercise 3.5.6 (Repeating lists) A list is **repeating** if it has an element appearing at two different positions. For instance, $[2, 5, 3, 5]$ is repeating and $[2, 5, 3]$ is not repeating.

- a) Declare a function testing whether a list is repeating.
- b) Declare a function testing whether a list is non-repeating.
- c) Declare a function that given a list l yields a non-repeating list containing the same elements as l .

3.6 Head and Tail

In OCaml we can declare polymorphic functions

$$hd : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \alpha$$

$$tl : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$$

that yield the head and the tail of nonempty lists:

```
let hd l =  
  match l with  
  | [] -> failwith "hd"  
  | x :: _ -> x  
  
let tl l =  
  match l with  
  | [] -> failwith "tl"  
  | _ :: l -> l
```

3 Lists

Both functions raise exceptions when applied to the empty list. Note the use of the underline symbol “_” for pattern variables that are not used in a rule. Also note the use of the predefined function

$$failwith : \forall \alpha. \text{string} \rightarrow \alpha$$

raising an exception *Failure s* carrying the string *s* given as argument. Interesting is the polymorphic type of *failwith* making it possible to type an application of *failwith* with whatever type is required by the context of the application.

At the mathematical level we don't admit exceptions. Thus we cannot define a polymorphic function

$$hd : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \alpha$$

since we don't have a value we can return for the empty list over α .

Exercise 3.6.1 Define a polymorphic function $tl : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$ returning the tail of nonempty lists. Do not use exceptions.

Exercise 3.6.2 Declare a polymorphic function $last : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \alpha$ returning the last element of a nonempty lists.

3.7 Position Lookup

The **positions** of a list are counted from left to right starting with the number 0. For instance, the list $[5, 6, 5]$ has the positions 0, 1, 2; moreover, the element at position 1 is 6, and the value 5 appears at the positions 0 and 2. The empty list has no position. More generally, a list of length n has the positions $0, \dots, n-1$.

We declare a tail-recursive *lookup function*

$$nth : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \text{int} \rightarrow \alpha$$

that given a list and a position returns the element at the position:

```
let rec nth l n =
  match l with
  | [] -> failwith "nth"
  | x :: l -> if n < 1 then x else nth l (n-1)
```

The function raises an exception if $l = []$ or $n \geq \text{length } l$. Note that $nth\ l\ 0$ yields the head of l if l is nonempty. Also note that nth terminates since it recurses on the list argument. Here is a trace:

$$nth\ [0, 2, 3, 5]\ 2 = nth\ [2, 3, 5]\ 1 = nth\ [3, 5]\ 0 = 3$$

3 Lists

Exercise 3.7.1 What is the result of $nth\ [0, 2, 3, 5]\ (-2)$?

Exercise 3.7.2 Declare a function $nth_checked$ that raises an invalid-argument exception (see §1.13) if $n < 0$ and otherwise agrees with nth . Check that your function behaves the same as the predefined function $List.nth$.

Exercise 3.7.3 Define a function $pos : \forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathbb{Z} \rightarrow \mathbb{B}$ testing whether a number is a position of a list.

Exercise 3.7.4 Declare a function $find : \forall\alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{N}$ that returns the first position of a list a given value appears at. For instance, we want $find\ 1\ [3, 1, 1] = 1$. If the value doesn't appear in the list, a failure exception should be raised.

3.8 Option Types

The lookup function for lists

$$nth : \forall\alpha. \mathcal{L}(\alpha) \rightarrow int \rightarrow \alpha$$

raises an exception if the position argument is not valid for the given list. There is the possibility to avoid the exception by changing the result type of nth to an option type

$$nth_opt : \forall\alpha. \mathcal{L}(\alpha) \rightarrow int \rightarrow \mathcal{O}(\alpha)$$

that in addition to the values of α has an extra value **None** that can be used to signal that the given position is not valid. In fact, OCaml comes with a type family

$$\mathcal{O}(\alpha) ::= \text{None} \mid \text{Some } \alpha$$

whose values are obtained with two polymorphic constructors

$$\text{Some} : \forall\alpha. \alpha \rightarrow \mathcal{O}(\alpha)$$

$$\text{None} : \forall\alpha. \mathcal{O}(\alpha)$$

such that **Some** injects the values of a type t into $\mathcal{O}(t)$ and **None** represents the extra value. We can declare a lookup function returning options as follows:

```
let rec nth_opt l n =
  match l with
  | [] -> None
  | x :: l -> if n < 1 then Some x else nth_opt l (n-1)
```

3 Lists

An option may be seen as a list that is either empty or a singleton list $[x]$. In fact, it is possible to replace option types with list types. Doing this gives away information as it comes to type checking.

Exercise 3.8.1 Declare a function *nth_opt_checked* that raises an invalid-argument exception if $n < 0$ and otherwise agrees with *nth_opt*.

Exercise 3.8.2 Declare a function

$$\text{nth_list} : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \text{int} \rightarrow \mathcal{L}(\alpha)$$

that agrees with *nth_opt* but returns a list with at most one element.

Exercise 3.8.3 Declare a function *find_opt* : $\forall \alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{O}(\mathbb{N})$ that returns the first position of a list a given value appears at. For instance, we want *find_opt* 7 [3, 7, 7] = *Some* 1 and *find_opt* 2 [3, 7, 7] = *None*.

3.9 Generalized Match Expressions

OCaml provides pattern matching in more general form than the basic list matches we have seen so far. A good example for explaining **generalized match expressions** is a function

$$\text{eq} : \forall \alpha. (\alpha \rightarrow \alpha \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B}$$

testing equality of lists using an equality test for the base type given as argument:

```
let rec eq (p: 'a -> 'a -> bool) l1 l2 =
  match l1, l2 with
  | [], [] -> true
  | x::l1, y::l2 -> p x y && eq p l1 l2
  | _, _ -> false
```

The generalized match expression in the declaration matches on two values and uses a final **catch-all rule**. Evaluation of a generalized match tries the patterns of the rules in the order they are given and commits to the first rule whose pattern matches. The pattern of the final rule will always match and will thus be used if no other rule applies. One speaks of a *catch all rule*.

Note that the above declaration specifies the type of the argument p using a type variable α . Without this specification OCaml would infer the more general type $\forall \alpha \beta. (\alpha \rightarrow \beta \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\beta) \rightarrow \mathbb{B}$ for *eq*.

The generalized match in the above declaration translates to a simple match with nested simple matches:

3 Lists

```
match l1 with
| [] ->
  begin match l2 with
    | [] -> true
    | _ :: _ -> false
  end
| x::l1 ->
  begin match l2 with
    | [] -> false
    | y::l2 -> p x y && eq p l1 l2
  end
end
```

The keywords **begin** and **end** provide a notational variant for a pair (\dots) of parentheses.

The notions of *disjointness* and *exhaustiveness* established for defining equations in §1.6 carry over to generalized match expressions. We will only use exhaustive match expressions but occasionally use non-disjoint match expressions where the order of the rules matters (e.g., catch-all rules). We remark that simple match expressions are always disjoint and exhaustive.

We see generalized match expressions as derived forms that compile into simple match expressions.

OCaml also has match expressions for tuples. For instance,

```
let fst a =
  match a with
  | (x, _) -> x
```

declares a projection function $\text{fst} : \forall \alpha \beta. \alpha \times \beta \rightarrow \alpha$ for pairs. In fact, match expressions for tuples are native in OCaml and the uses of tuple patterns we have seen in let expressions, lambda abstractions, and declarations all compile into match expressions for tuples.

Patterns in OCaml may also contain numbers and other constants. For instance, we may declare a function that tests whether a list starts with the numbers 1 and 2 as follows:

```
let test l =
  match l with
  | 1 :: 2 :: _ -> true
  | _ -> false
```

Exercise 3.9.1 Declare a function testing whether a list starts with the numbers 1 and 2 just using simple match expressions for lists.

Exercise 3.9.2 Declare a function $\text{swap} : \forall \alpha \beta. \alpha \times \beta \rightarrow \beta \times \alpha$ swapping the components of a pair using a simple match expression for tuples.

Exercise 3.9.3 (Maximal element) Declare a function that yields the maximal element of a list of numbers. If the list is empty, a failure exception should be raised.

Exercise 3.9.4 Translate the expression

```
fun l -> match l with
| 0::x::_ -> Some x
| x::1::_ -> Some x
| _ -> None
```

into an expression only using simple matches.

3.10 Sublists

A **sublist** of a list l is obtained by deleting $n \geq 0$ positions of l . For instance, the sublists of $[1, 2]$ are the lists

$[1, 2], [2], [1], []$

We observe that the empty list $[]$ has only itself as a sublist, and that a sublist of a nonempty list $x :: l$ is either a sublist of l , or a list $x :: l'$ where l' is a sublist of l . Using this observation, it is straightforward to define a function that yields a list of all sublists of a list:

$$\begin{aligned} pow &: \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\mathcal{L}(\alpha)) \\ pow [] &:= [[]] \\ pow (x :: l) &:= pow\ l\ @\ map\ (\lambda l. x :: l)\ (pow\ l) \end{aligned}$$

We call $pow\ l$ the **power list** of l .

A sublist test “ l_1 is sublist of l_2 ” needs a more involved case analysis:

$$\begin{aligned} is_sublist &: \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B} \\ is_sublist\ l_1\ [] &:= l_1 = [] \\ is_sublist\ []\ (y :: l_2) &:= \text{true} \\ is_sublist\ (x :: l_1)\ (y :: l_2) &:= is_sublist\ (x :: l_1)\ l_2\ || \\ &\quad x = y \ \&\& \ is_sublist\ l_1\ l_2 \end{aligned}$$

We remark that a computationally naive sublist test can be obtained with the power list function and the membership test:

$$is_sublist\ l_1\ l_2 = mem\ l_1\ (pow\ l_2)$$

Exercise 3.10.1 (Graded power list)

Declare a function $gpow : \forall \alpha. \mathbb{N} \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\mathcal{L}(\alpha))$ such that $gpow\ k\ l$ yields a list containing all sublists of l of length k .

3 Lists

Exercise 3.10.2 (Prefixes, Segments, Suffixes)

Given a list $l = l_1 @ l_2 @ l_3$, we call l_1 a **prefix**, l_2 a **segment**, and l_3 a **suffix** of l . The definition is such that prefixes are segments starting at the beginning of a list, and suffixes are segments ending at the end of a list. Moreover, every list is a prefix, segment, and suffix of itself.

- Convince yourself that segments are sublists.
- Give a list and a sublist that is not a segment of the list.
- Declare a function that yields a list containing all prefixes of a list.
- Declare a function that yields a list containing all suffixes of a list.
- Declare a function that yields a list containing all segments of a list.

Exercise 3.10.3 (Splits) Given a list $l = l_1 @ l_2$, we call the pair (l_1, l_2) a **split** of l . Declare a function that yields a list containing all splits of a list.

Exercise 3.10.4 Declare a function $filter : \forall \alpha. (\alpha \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$ that given a test and a list yields the sublist of all elements that pass the test. For instance, we want $filter (\lambda x. x > 2) [2, 5, 1, 5, 2] = [5, 5]$.

3.11 Folding Lists

Consider the list

$$a_1 :: (a_2 :: (a_3 :: []))$$

If we replace cons with $+$ and nil with 0, we obtain the expression

$$a_1 + (a_2 + (a_3 + 0))$$

which evaluates to the sum of the elements of the list. If we replace cons with \cdot and nil with 1, we obtain the expression

$$a_1 \cdot (a_2 \cdot (a_3 \cdot 1))$$

which evaluates to the product of the elements of the list. More generally, we obtain the expression

$$f a_1 (f a_2 (f a_3 b))$$

if we replace cons with a function f and nil with a value b . Even more generally, we can define a function $fold$ such that $fold f l b$ yields the value of the expression obtained from l by replacing cons with f and nil with b :

$$\begin{aligned} fold &: \forall \alpha \beta. (\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \mathcal{L}(\alpha) \rightarrow \beta \rightarrow \beta \\ fold f [] b &:= b \\ fold f (a :: l) b &:= f a (fold f l b) \end{aligned}$$

3 Lists

We have the following equations:

$$\begin{aligned}
\text{fold } (+) [x_1, \dots, x_n] 0 &= x_1 + \dots + x_n + 0 \\
\text{fold } (\lambda ab. a^2 + b) [x_1, \dots, x_n] 0 &= x_1^2 + \dots + x_n^2 + 0 \\
l_1 @ l_2 &= \text{fold } (::) l_1 l_2 \\
\text{flatten } l &= \text{fold } (@) l [] \\
\text{length } l &= \text{fold } (\lambda ab. b + 1) l 0 \\
\text{rev } l &= \text{fold } (\lambda ab. b @ [a]) l []
\end{aligned}$$

Folding of lists is similar to iteration with numbers in that both recursion schemes can express many functions without further recursion.

There is a tail-recursive variant *foldl* of *fold* satisfying the equation $\text{foldl } f \ l \ b = \text{fold } f \ (\text{rev } l) \ b$:

$$\begin{aligned}
\text{foldl} : \forall \alpha \beta. (\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \mathcal{L}(\alpha) \rightarrow \beta \rightarrow \beta \\
\text{foldl } f \ [] \ b &:= b \\
\text{foldl } f \ (a :: l) \ b &:= \text{foldl } f \ l \ (f \ a \ b)
\end{aligned}$$

One says that *fold* folds a list from the right

$$\text{fold } f \ [a_1, a_2, a_3] \ b = f \ a_1 \ (f \ a_2 \ (f \ a_3 \ b))$$

that *foldl* folds a list from the left

$$\text{foldl } f \ [a_1, a_2, a_3] \ b = f \ a_3 \ (f \ a_2 \ (f \ a_1 \ b))$$

If the order of the folding is not relevant, the tail-recursive version *foldl* is preferable over *fold*.

OCaml provides the function *fold* as *List.fold_right*. OCaml also provides a function *List.fold_left*, which however varies the argument order of our function *foldl*. To avoid confusion, we will not use *List.fold_left* in this chapter but instead use our function *foldl*.

Exercise 3.11.1 Using *fold*, declare functions that yield the concatenation, the flattening, the length, and the reversal of lists.

Exercise 3.11.2 Using *foldl*, declare functions that yield the length, the reversal, and the concatenation of lists.

Exercise 3.11.3 We have the equations

$$\begin{aligned}
\text{foldl } f \ l \ b &= \text{fold } f \ (\text{rev } l) \ b \\
\text{fold } f \ l \ b &= \text{foldl } f \ (\text{rev } l) \ b
\end{aligned}$$

- Show that the second equation follows from the first equation using the equation $\text{rev}(\text{rev } l) = l$.
- Obtain *fold* from *foldl* not using recursion.
- Obtain *foldl* from *fold* not using recursion.

3.12 Insertion Sort

A sequence x_1, \dots, x_n of numbers is called **sorted** if its elements appear in order: $x_1 \leq \dots \leq x_n$. **Sorting** a sequence means to rearrange the elements such that the sequence becomes sorted. We want to define a function

$$\text{sort} : \mathcal{L}(\mathbb{Z}) \rightarrow \mathcal{L}(\mathbb{Z})$$

such that $\text{sort } l$ is a sorted rearrangement of l . For instance, we want $\text{sort } [5, 3, 2, 7, 2] = [2, 2, 3, 5, 7]$.

For now, we only consider sorting for lists of numbers. Later it will be easy to generalize to other types and other orders.

There are different sorting algorithms. Probably the easiest one is **insertion sort**. For insertion sort one first defines a function that inserts a number x into a list such that the result list is sorted if the argument list is sorted.² Now sorting a list l is easy: We start with the empty list, which is sorted, and insert the elements of l one by one. Once all elements are inserted, we have a sorted rearrangement of l . Here are definitions of the necessary functions:

$$\text{insert} : \mathbb{Z} \rightarrow \mathcal{L}(\mathbb{Z}) \rightarrow \mathcal{L}(\mathbb{Z})$$

$$\text{insert } x [] := [x]$$

$$\text{insert } x (y :: l) := \text{IF } x \leq y \text{ THEN } x :: y :: l \text{ ELSE } y :: \text{insert } x l$$

$$\text{isort} : \mathcal{L}(\mathbb{Z}) \rightarrow \mathcal{L}(\mathbb{Z})$$

$$\text{isort} [] := []$$

$$\text{isort } (x :: l) := \text{insert } x (\text{isort } l)$$

Make sure you understand every detail of the definition. We offer a trace:

$$\begin{aligned} \text{isort } [3, 2] &= \text{insert } 3 (\text{isort } [2]) \\ &= \text{insert } 3 (\text{insert } 2 (\text{isort } [])) \\ &= \text{insert } 3 (\text{insert } 2 []) \\ &= \text{insert } 3 [2] = 2 :: \text{insert } 3 [] = 2 :: [3] = [2, 3] \end{aligned}$$

Note that isort inserts the elements of the input list reversing the order they appear in the input list.

²More elegantly, we may say that the insertion function preserves sortedness.

3 Lists

Comparisons are polymorphically typed

Declaring the functions *insert* and *isort* in OCaml is now routine. There is, however, the surprise that OCaml derives polymorphic types for *insert* and *isort* if no type specification is given:

$$\begin{aligned} \textit{insert} &: \forall\alpha. \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \\ \textit{isort} &: \forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \end{aligned}$$

This follows from the fact that OCaml accommodates comparisons with the polymorphic type

$$\forall\alpha. \alpha \rightarrow \alpha \rightarrow \mathbb{B}$$

it also uses for the equality test (§2.6). We will explain later how comparisons behave on tuples and lists. For boolean values, OCaml realizes the order $\text{false} < \text{true}$. Thus we have

$$\textit{isort} [\text{true}, \text{false}, \text{false}, \text{true}] = [\text{false}, \text{false}, \text{true}, \text{true}]$$

Exercise 3.12.1 Declare a function *sorted* : $\forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathbb{B}$ that tests whether a list is sorted. Use tail recursion. Write the function with a generalized match and show how the generalized match translates into simple matches.

Exercise 3.12.2 Declare a function *perm* : $\forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B}$ that tests whether two lists are equal up to reordering.

Exercise 3.12.3 (Sorting into descending order)

Declare a function *sort_desc* : $\forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$ that reorders a list such that the elements appear in descending order. For instance, we want *sort_desc* [5, 3, 2, 5, 2, 3] = [5, 5, 3, 3, 2, 2].

Exercise 3.12.4 (Sorting with duplicate deletion)

Declare a function *dsort* : $\forall\alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$ that sorts a list and removes all duplicates. For instance, *dsort* [5, 3, 2, 5, 2, 3] = [2, 3, 5].

Insertion order

Sorting by insertion inserts the elements of the input list one by one into the empty list. The order in which this is done does not matter for the result. The function *isort* defined above inserts the elements of the input list reversing the order of the input list. If we define *isort* as

$$\textit{isort} \, l := \textit{fold} \, \textit{insert} \, l \, []$$

3 Lists

we preserve the insertion order. If we switch to the definition

$$isort\ l := foldl\ insert\ l\ []$$

we obtain a tail-recursive insertion function inserting the elements of the input list in the order they appear in the input list.

Exercise 3.12.5 (Count Tables) Declare a function

$$table : \forall \alpha. \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha \times \mathbb{N}^+)$$

such that $(x, n) \in table\ l$ if and only if x occurs $n > 0$ times in l . For instance, we want

$$table\ [4, 2, 3, 2, 4, 4] = [(4, 3), (2, 2), (3, 1)]$$

Make sure *table* lists the count pairs for the elements of l in the order the elements appear in l , as in the example above.

3.13 Generalized Insertion Sort

Rather than sorting lists using the predefined order \leq , we may sort lists using an order given as argument:

$$insert : \forall \alpha. (\alpha \rightarrow \alpha \rightarrow \mathbb{B}) \rightarrow \alpha \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$$

$$insert\ p\ x\ [] := [x]$$

$$insert\ p\ x\ (y :: l) := \text{IF } p\ x\ y \text{ THEN } x :: y :: l \text{ ELSE } y :: insert\ p\ x\ l$$

$$gisort : \forall \alpha. (\alpha \rightarrow \alpha \rightarrow \mathbb{B}) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha)$$

$$gisort\ p\ l := fold\ (insert\ p)\ l\ []$$

Now the function *gisort* (\leq) sorts as before in ascending order, while the function *gisort* (\geq) sorts in descending order:

$$gisort\ (\geq)\ [1, 3, 3, 2, 4] = [4, 3, 3, 2, 1]$$

When we declare the functions *insert* and *gisort* in OCaml, we can follow the mathematical definitions. Alternatively, we can declare *gisort* using a local declaration for *insert*:

```
let gisort p l =
  let rec insert x l =
    match l with
    | [] -> [x]
    | y :: l -> if p x y then x :: y :: l else y :: insert x l
  in
  foldl insert l []
```

3 Lists

This way we avoid the forwarding of the argument p .

Exercise 3.13.1 Declare a function

$$\text{reorder} : \forall \alpha \beta. \mathcal{L}(\alpha \times \beta) \rightarrow \mathcal{L}(\alpha \times \beta)$$

that reorders a list of pairs such that the first components of the pairs are ascending. If there are several pairs with the same first component, the original order of the pairs should be preserved. For instance, we want $\text{reorder} [(5, 3), (3, 7), (5, 2), (3, 2)] = [(3, 7), (3, 2), (5, 3), (5, 2)]$. Declare reorder as a one-liner using the sorting function gisort .

3.14 Lexical Order

We now explain how we obtain an order for lists over t from an order for the base type t following the principle used for ordering words in dictionaries. We speak of a **lexical ordering**. Examples for the lexical ordering of lists of integers are

$$[] < [-1] < [-1, -2] < [0] < [0, 0] < [0, 1] < [1]$$

The general principle behind the lexical ordering can be formulated with two rules:

- $[] < x :: l$
- $x_1 :: l_1 < x_2 :: l_2$ if either $x_1 < x_2$, or $x_1 = x_2$ and $l_1 < l_2$.

Following the rules, we define a function that yields a test for the lexical order \leq of lists given a test for an order \leq of the base type:

$$\begin{aligned} \text{lex} : \forall \alpha. (\alpha \rightarrow \alpha \rightarrow \mathbb{B}) &\rightarrow \mathcal{L}(\alpha) \rightarrow \mathcal{L}(\alpha) \rightarrow \mathbb{B} \\ \text{lex } p [] l_2 &:= \text{true} \\ \text{lex } p (x_1 :: l_1) [] &:= \text{false} \\ \text{lex } p (x_1 :: l_1) (x_2 :: l_2) &:= p \ x_1 \ x_2 \ \&\& \\ &\quad \text{IF } p \ x_2 \ x_1 \ \text{THEN } \text{lex } p \ l_1 \ l_2 \ \text{ELSE } \text{true} \end{aligned}$$

Exercise 3.14.1 (Lexical order for pairs) The idea of lexical order extends to pairs and to tuples in general.

- Explain the lexical order of pairs of type $t_1 \times t_2$ given orders for the component types t_1 and t_2 .
- Declare a function

$$\text{lexP} : \forall \alpha \beta. (\alpha \rightarrow \alpha \rightarrow \mathbb{B}) \rightarrow (\beta \rightarrow \beta \rightarrow \mathbb{B}) \rightarrow \alpha \times \beta \rightarrow \alpha \times \beta \rightarrow \mathbb{B}$$

testing the lexical order of pairs. For instance, we want

$$\text{lexP } (\leq) (\geq) (1, 2) (1, 3) = \text{false}$$

and $\text{lexP } (\leq) (\geq) (0, 2) (1, 3) = \text{true}$.

3.15 Prime Factorization

Every integer greater than 1 can be written as a product of prime numbers; for instance,

$$\begin{aligned} 60 &= 2 \cdot 2 \cdot 3 \cdot 5 \\ 147 &= 3 \cdot 7 \cdot 7 \\ 735 &= 3 \cdot 5 \cdot 7 \cdot 7 \end{aligned}$$

One speaks of the *prime factorization* of a number. Recall that a prime number is an integer greater than 1 that cannot be obtained as the product of two integers greater than 1 (§2.4).

It is straightforward to compute the smallest prime factor of a number $x \geq 2$: We simply search for the first $k \geq 2$ **dividing** x (i.e., $x \% k = 0$). Such a k exists since x divides x . Moreover, the first k dividing x is always prime since otherwise there would be a smaller number greater than 1 dividing k and thus also x .

If we can compute smallest prime factors, we can compute prime factorizations by dividing with the prime factor found and continuing recursively. We realize this algorithm with an OCaml function

$$\text{prime_fac} : \text{int} \rightarrow \mathcal{L}(\text{int})$$

declared as follows:

```
let rec prime_fac x =
  if x < 2 then []
  else let k = first (fun k -> x mod k = 0) 2 in
    if k = x then [x]
    else k :: prime_fac (x / k)
```

We have $\text{prime_fac } 735 = [3, 5, 7, 7]$, for instance.

As is, the algorithm is slow on large prime numbers (try 479,001,599). The algorithm can be made much faster by stopping the linear search for the first k dividing x once $k^2 > x$ since then the first k dividing x is x . Moreover, if we have a least prime factor $k < x$, it suffices to start the search for the next prime factor at k since we know that no number smaller than k divides x . We realize the optimized algorithm as follows:

```
let rec prime_fac' k x =
  if k * k > x then [x]
  else if x mod k = 0 then k :: prime_fac' k (x / k)
  else prime_fac' (k + 1) x
```

3 Lists

We have $\text{prime_fac}'\ 2\ 735 = [3, 5, 7, 7]$, for instance. Interestingly, the optimized algorithm is simpler than the naive algorithm we started from (as it comes to code, but not as it comes to correctness). It makes sense to define a wrapper function for $\text{prime_fac}'$ ensuring that $\text{prime_fac}'$ is applied with admissible arguments:

```
let prime_fac x = if x < 2 then [] else prime_fac' 2 x
```

We used several mathematical facts to derive the optimized prime factorization algorithm:

1. If $2 \leq x < k^2$ and no number $2 \leq n \leq k$ divides x , then x is a prime number.
2. If $2 \leq k < x$, and k divides x , and no number $2 \leq n < k$ divides x , then k is the least prime factor of x .
3. If $2 \leq k < x$, and k divides x , and no number $2 \leq n < k$ divides x , then no number $2 \leq n < k$ divides x/k .

The correctness of the algorithm also relies on the fact that the **safety condition**

- $2 \leq k \leq x$ and no number $2 \leq n < k$ divides x

propagates from every initial application of $\text{prime_fac}'$ to all recursive applications. We say that the safety condition is an **invariant** for the applications of $\text{prime_fac}'$.

It suffices to argue the termination of $\text{prime_fac}'$ for the case that the safety condition is satisfied. In this case the $x - k \geq 0$ is decreased by every recursion step.

Exercise 3.15.1 Give traces for the following applications:

- a) $\text{prime_fac}'\ 2\ 7$ b) $\text{prime_fac}'\ 2\ 8$ c) $\text{prime_fac}'\ 2\ 15$

Exercise 3.15.2 Declare a function that yields the least prime factor of an integer $x \geq 2$. Make sure that at most $\sqrt[3]{x}$ remainder operations are necessary.

Exercise 3.15.3 Declare a primality test using at most $\sqrt[3]{x}$ remainder operations for an argument $x \geq 2$.

Exercise 3.15.4 Dieter Schlau has simplified the naive prime factorization function:

```
let rec prime_fac x =
  if x < 2 then []
  else let k = first (fun k -> x mod k = 0) 2 in
    k :: prime_fac (x / k)
```

Explain why Dieter's function is correct.

3.16 Key-Value Maps

Recall the use of environments in the typing and evaluation systems for expressions we discussed in Chapter ???. There environments are modeled as sequences of pairs binding variables to types or values, as for instance in the value environment $[x \triangleright 5, y \triangleright 7, z \triangleright 2]$. Abstractly, we can see an environment as a list of pairs

$$(key, value)$$

consisting of a *key* and a *value*, where the key may be a string or a number representing a variable, and the value may represent a type or an evaluation value. One says that key-value lists are **maps** from keys to values. Following this consideration, we now define **maps** as values of the type family

$$map\ \alpha\ \beta := \mathcal{L}(\alpha \times \beta)$$

Recall that the typing and evaluation rules need only two operations

$$\begin{aligned} lookup &: \forall \alpha \beta. map\ \alpha\ \beta \rightarrow \alpha \rightarrow \mathcal{O}(\beta) \\ update &: \forall \alpha \beta. map\ \alpha\ \beta \rightarrow \alpha \rightarrow \beta \rightarrow map\ \alpha\ \beta \end{aligned}$$

on environments, where *lookup* yields the value for a given key provided the environment contains a pair for the key, and *update* updates an environment with a given key-value pair. For instance,

$$\begin{aligned} lookup\ [("x", 5), ("y", 13), ("z", 2)]\ "y" &= 13 \\ update\ [("x", 5), ("y", 7), ("z", 2)]\ "y"\ 13 &= [("x", 5), ("y", 13), ("z", 2)] \end{aligned}$$

The defining equations for *lookup* and *update* are as follows:

$$\begin{aligned} lookup\ []\ a &:= \text{None} \\ lookup\ ((a', b) :: l)\ a &:= \text{IF } a' = a \text{ THEN Some } b \\ &\quad \text{ELSE } lookup\ l\ a \\ update\ []\ a\ b &:= [(a, b)] \\ update\ ((a', b') :: l)\ a\ b &:= \text{IF } a' = a \text{ THEN } (a, b) :: l \\ &\quad \text{ELSE } (a', b') :: update\ l\ a\ b \end{aligned}$$

Exercise 3.16.1 Give the values of the following expressions:

- a) $update\ (update\ (update\ []\ "x"\ 7)\ "y"\ 2)\ "z"\ 5$
- b) $lookup\ (update\ l\ "x"\ 13)\ "x"$
- c) $lookup\ (update\ l\ a\ 7)\ a$

3 Lists

Exercise 3.16.2 Decide for each of the following equations whether it is true in general.

- a) $\text{lookup } (\text{update } l \ a \ b) \ a = \text{Some } b$
- b) $\text{lookup } (\text{update } l \ a' \ b) \ a = \text{lookup } l \ a \quad \text{if } a' \neq a$
- c) $\text{update } (\text{update } l \ a \ b) \ a' \ b' = \text{update } (\text{update } l \ a' \ b') \ a \ b \quad \text{if } a \neq a'$
- d) $\text{lookup } (\text{update } (\text{update } l \ a \ b) \ a' \ b') \ a = \text{Some } b \quad \text{if } a \neq a'$

Exercise 3.16.3 (Boundedness) Declare a function

$$\text{bound} : \forall \alpha \beta. \text{map } \alpha \ \beta \rightarrow \alpha \rightarrow \text{bool}$$

that checks whether a map binds a given key. Note that you can define *bound* using *lookup*.

Exercise 3.16.4 (Deletion) Declare a function

$$\text{delete} : \forall \alpha \beta. \text{map } \alpha \ \beta \rightarrow \alpha \rightarrow \text{map } \alpha \ \beta$$

deleting the entry for a given key. We want $\text{lookup } (\text{delete } l \ a) \ a = \text{None}$ for all environments l and all keys a .

Exercise 3.16.5 (Maps with memory) Note that *lookup* searches maps from left to right until it finds a pair with the given key. This opens up the possibility to keep previous values in the map by modifying *update* so that it simply appends the new key-value pair in front of the list:

$$\text{update } l \ a \ b := (a, b) :: l$$

Redo the previous exercises for the new definition of *update*. Also define a function

$$\text{lookup_all} : \forall \alpha \beta. \text{map } \alpha \ \beta \rightarrow \alpha \rightarrow \mathcal{L}(\beta)$$

that yields the list of all values for a given key.

3 Lists

Exercise 3.16.6 (Maps as functions) Maps can be realized with functions if all maps are constructed from the empty map

$$\text{empty} : \forall \alpha \beta. \text{map } \alpha \beta$$

with *update* and the only thing that matters is that *lookup* yields the correct results. Assume the definition

$$\text{map } \alpha \beta := \alpha \rightarrow \text{Some } \beta$$

and define *empty* and the operations *update* and *lookup* accordingly. Test your solution with

$$\text{lookup } (\text{update } (\text{update } (\text{update } \text{empty } \text{"y"} 7) \text{"x"} 2) \text{"y"} 5) \text{"y"} = \text{Some } 5$$

Note that you can still define the operations *bound* and *delete* from Exercises 3.16.1 and 3.16.2.