
Programming in homotopy type theory and erasing propositions

Gabe Dijkstra

M.Sc. thesis ICA-3354881
[Supervisors] Wouter Swierstra and Johan Jeuring

August 20, 2013



Universiteit Utrecht

Department of Computing Science

Abstract

In the recent years, homotopy type theory has become the subject of much study. It studies the correspondence between propositional equality of Martin-Löf's type theory and the concept of homotopy from topology. This roughly means that inhabitants of a type can be seen as points of a space and that a propositional equality $x \equiv y$ can be seen as a path $x \rightsquigarrow y$. This thesis aims to provide an introduction to homotopy type theory geared toward programmers familiar with dependently typed programming, but unfamiliar with topology. We will present applications of homotopy type theory to programming, such as quotient types and dealing with views on abstract types. Apart from this, we will discuss the use of h -propositions to identify parts of a program that are not needed at runtime and discuss how this can be used to optimise our programs.

Contents

1	Introduction	2
2	Homotopy type theory	4
2.1	Homotopy theory	4
2.2	Identity types of Martin-Löf's type theory	5
2.3	Homotopy interpretation	8
2.4	n -types and truncations	10
2.5	Higher inductive types	12
2.6	Equivalence and univalence	16
2.7	Implementation	17
3	Applications of homotopy type theory	19
3.1	Quotient types	19
3.2	Views on abstract types	24
3.3	Conclusion	30
4	Erasing propositions	31
4.1	Propositions	31
4.2	The <i>Prop</i> universe in Coq	33
4.3	Irrelevance in Agda	36
4.4	Collapsible families	38
4.5	Internalising collapsibility	40
4.6	Internalising the collapsibility optimisation	42
4.7	Indexed h -propositions and homotopy type theory	44
4.8	Conclusions	46
5	Discussion	48
	Index of symbols	52
	Index	53

Chapter 1

Introduction

One of the tricky things that comes up when designing a type system or a logic, is the defining a right notion of equality. When type checking a term, one needs a notion of equality, called *definitional equality*, in this thesis denoted by \triangleq . For example when one type checks an application $f\ a$ and we know that $f : A \rightarrow B$ and we know that $a : X$, we have to check that A and X are equal in some way. In Martin-Löf's type theory, A and X need to be definitionally equal: if we reduce both A and X to their normal forms, they need to be syntactically equal.

We also want to be able to reason about equality in the type theory itself, e.g. use it to show that two programs behave in the same way, when given the same input. The notion of equality internal to the type theory is called *propositional equality* (in this thesis denoted by \equiv). In Martin-Löf's type theory, propositional equality is defined using the so called *identity types*: an inductive family with *refl* as its only constructor. This construction essentially imports definitional equality into the type theory.

However, the resulting structure is not exactly definitional equality. We can force the two notions to coincide by adding an *equality reflection* rule, i.e. a rule that states that if we have a proof that x and y are propositionally equal, they are also definitionally equal. Since type checking makes use of definitional equality, to show that two terms are definitionally equal, we may need to produce a proof of propositional equality first. This proof search means that type checking becomes undecidable. However, adding equality reflection does mean that we can prove useful things such as function extensionality $((x : A) \rightarrow f\ x \equiv g\ x) \rightarrow f \equiv g$, something that we cannot prove if we leave the equality reflection rule out.

The study of intensional type theory, i.e. type theory without the equality reflection rule, involved finding out why we cannot prove certain properties about propositional equality that were deemed to be natural properties for a notion of equality. This eventually led to the discovery of homotopy type theory, an interpretation of types and their identity types in the language of homotopy theory:

type theory	homotopy theory
A is a type	A is a space
$x, y : A$	x and y are points in A
$p, q : x \equiv y$	p and q are paths from x to y
$w : p \equiv q$	w is a homotopy between paths p and q

The discovery was that propositional equality behaves just like the homotopy we know from topology. This discovery spawned a lot of interest, as it meant that the language of type theory can be used to prove theorems about homotopy theory. It also means that we can use our intuition about homotopy theory to make statements about type theory. For example, one question that has remained unanswered for nearly two decades, the question why uniqueness of identity proofs cannot be proven using J , can now be answered by drawing a couple of pictures (see subsection 2.3.1).

This leads us to the main research question:

What is homotopy type theory and why is it interesting to do programming in?

In chapter 2 we give an introduction and overview of some of the main concepts of homotopy type theory. In chapter 3 we discuss several applications of homotopy type theory to programming. In particular we look at how we can implement quotient types in homotopy type theory and contrast this to popular ways to work with quotient types. Another application we consider is the use of univalence to deal with views on abstract types. We work out the example given by [Licata \[2012\]](#) and extend the result to non-isomorphic views, using quotient types.

Homotopy type theory provides us with a notion of propositions, the so called h -propositions. In chapter 4 we compare this to similar notions found in Coq, Agda and Epigram.

We investigate whether we can formulate an optimisation based on h -propositions in the spirit of the collapsibility optimisation proposed in [Brady et al. \[2004\]](#).

In the final chapter, chapter 5, we will discuss our answers to our research questions and propose directions of future research.

Since the focus of this thesis is on the programming aspects of homotopy type theory, as opposed to doing homotopy theory, we will not do any diagram chasing and instead will use Agda syntax throughout the thesis. As such, we will expect the reader to be familiar with this language.

Notation The code in the thesis will not always be valid Agda code. We will use the notation $A : Type$ instead of $A : Set$, in order to avoid confusion between types and the homotopy type theory notion of h -sets. We will also refrain from mentioning levels and essentially assume $Type : Type$. The accompanying will explicitly mention the levels.

Code Look at `README.agda`.

Chapter 2

Homotopy type theory

As was briefly mentioned in chapter 1, homotopy type theory studies the correspondence between homotopy theory and type theory. As such, we will start out with a very brief sketch of the basic notions of homotopy theory (section 2.1). After that, we will describe the notion of propositional equality in Martin-Löf's type theory using identity types (section 2.2). Having defined the identity types, we can explain the interpretation of Martin-Löf's type theory in homotopy theoretic terms, relating propositional equality to paths (section 2.3). In section 2.4 we describe how the idea of classifying spaces along their homotopic structure can be used in type theory to classify types. section 2.5 and section 2.6 describe two extensions to Martin-Löf's type theory inspired by homotopy theory. This chapter is concluded by a discussion on the implementation issues of homotopy type theory (section 2.7).

2.1 Homotopy theory

In *homotopy theory* we are interested in studying *continuous deformations*. The simplest case of this is continuously deforming one point into another point, which is called a *path*. A path in a space X from point x to y is a continuous function $p : [0, 1] \rightarrow X$, such that $p\ 0 = x$ and $p\ 1 = y$, also notated as $p : x \rightsquigarrow y$. The set of all paths in X can be also considered as a space. In this space, called the *path space* of X , we again can look at the paths. Suppose we have two paths $p, q : [0, 1] \rightarrow X$ with the same begin and end points, then a path between p and q , called a *homotopy*, is a continuous function $\gamma : [0, 1] \rightarrow [0, 1] \rightarrow X$ where $\gamma\ 0 = p$ and $\gamma\ 1 = q$ (see ??). Of course, we can also look at homotopies in these path spaces, and homotopies between these higher homotopies, ad infinitum.

If we have a path $p : a \rightsquigarrow b$ and a path $q : b \rightsquigarrow c$, we can compose these to form a path $p \circ q : a \rightsquigarrow c$. For every path $p : a \rightsquigarrow b$, there is a reversed path $p^{-1} : b \rightsquigarrow a$. For every point a , there is the constant path $r_a : a \rightsquigarrow a$. One might wonder whether reversing a path acts as an inverse operation with r_a being the unit of path composition, i.e. whether the following equations are satisfied:

- $p \circ p^{-1} = r_a$

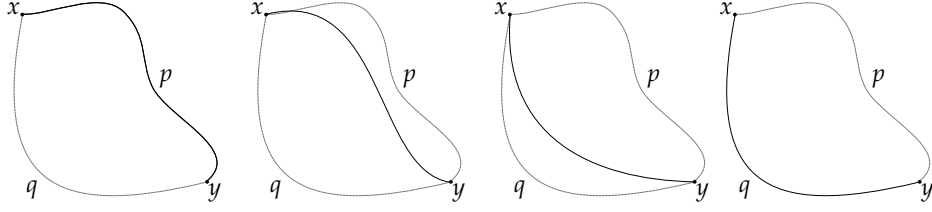


Figure 2.1: A homotopy between paths p and q

- $p^{-1} \circ p = r_b$
- $p \circ r_b = p$
- $r_a \circ p = p$

This happens to not be the case: the equations do not hold in the strict sense. However, both sides of the equations are homotopic to each other. The same holds for the associativity of composition: it is only associative up to homotopy. Homotopies also have this same structure, in which all the equalities hold up to higher homotopy. The groupoid-like structure that towers of homotopies have, is called a ∞ -groupoid structure. It was proposed by Grothendieck that homotopy theory should be the study of these ∞ -groupoids.

citation needed

2.2 Identity types of Martin-Löf's type theory

Martin-Löf [1985] introduced a notion of equality in his type theory: *propositional equality*, defined using so called identity types. These types can be formulated in Agda syntax as follows:

```
data Id (A : Type) : A → A → Type where
  refl : (x : A) → Id A x x
```

In order to type check $\text{refl } x : \text{Id } A \ x \ x$, the type checker needs to verify that x and y are definitionally equal. The refl constructor gives us that definitional equality implies propositional equality. The converse does not hold: we are working with an *intensional* type theory. *Extensional* type theories add a so called equality reflection rule that propositional equality implies definitional equality.

If we want to do something with the inhabitants of an inductive type, other than passing them around, we must use the induction principle (or elimination operator) of the inductive type. The induction principle of the Id type is usually called J and has the following type:

```
J : (A : Type)
  → (P : (x y : A) → (p : Id A x y) → Type)
  → (c : (x : A) → P x x (refl x))
  → (x y : A) → (p : Id A x y)
  → P x y p
```

Along with this type, we have the following computation rule:

$$J \ A \ P \ c \ x \ x \ (refl \ x) \triangleq c \ x$$

We will make use of a slightly different, but equivalent formulation of these types, due to Paulin-Mohring, where the x is a parameter as opposed to an index, yielding a more convenient elimination principle:

citation needed

data $Id' \ (A : Type) \ (x : A) : A \rightarrow Type$ **where**
 $refl : Id' \ A \ x \ x$

with induction principle:

$J' : (A : Type)$
 $\rightarrow (x : A)$
 $\rightarrow (P : (y : A) \rightarrow (p : Id' \ A \ x \ y) \rightarrow Type)$
 $\rightarrow (c : P \ x \ x \ refl)$
 $\rightarrow (y : A) \rightarrow (p : Id' \ A \ x \ y)$
 $\rightarrow P \ x \ y \ p$

and computation rule:

$$J' \ A \ P \ c \ x \ refl \triangleq c$$

To make things look more like the equations we are used to, we will for the most part use infix notation, leaving the type parameter implicit: $Id \ A \ x \ y$ becomes $x \equiv y$. In some cases we will fall back to the $Id \ A \ x \ y$ notation, when it is a bit harder to infer the type parameter.

Using the identity types and their induction principles, we can show that it is an equivalence relation, i.e. given $A : Type$ and $x \ y \ z : A$, we can find inhabitants of the following types:

- $refl : Id \ A \ x \ x$
- $symm : Id \ A \ x \ y \rightarrow Id \ A \ y \ x$
- $trans : Id \ A \ x \ y \rightarrow Id \ A \ y \ z \rightarrow Id \ A \ x \ z$

Another important property of propositional equality is that it is a congruence relation, i.e. we have a term with the following type:

$$ap : \{A \ B : Type\} \rightarrow (f : A \rightarrow B) \rightarrow \{x \ y : A\} \rightarrow x \equiv y \rightarrow f \ x \equiv f \ y$$

$ap \ f$ can be read as the (functorial) *action* on *paths* induced by f or the *application* of f on *paths*. If we want to generalise ap to also work on dependent functions $f : (a : A) \rightarrow B \ a$, we notice that we get something that does not type check: $f \ x \equiv f \ y$ does not type check because $f \ x : B \ x$ and $f \ y : B \ y$. However, if we have an equality between x and y , then $B \ x \equiv B \ y$, so we should be able to somehow transform something of type $B \ x$ to something of type $B \ y$. This process is called *transporting*:

$$transport : \{A : Type\} \{B : A \rightarrow Type\} \{x \ y : A\} \rightarrow x \equiv y \rightarrow B \ x \rightarrow B \ y$$

transport is sometimes also called *subst*, as *transport* witnesses the fact that if we have $x \equiv y$, we can substitute any occurrence of x in context B with y .

Using *transport* we can now formulate the dependent version of *ap*:

$$\begin{aligned} \text{apd} : \{ A : \text{Type} \} \{ B : A \rightarrow \text{Type} \} \{ x \ y : A \} \rightarrow (f : (a : A) \rightarrow B \ a) \rightarrow (\beta : x \equiv y) \\ \rightarrow \text{transport } \beta (f \ x) \equiv f \ y \end{aligned}$$

The resulting equality is an equality of between points in $B \ y$.

2.2.1 Difficulties of identity types

Even though at first glance the identity types have the right structure: they form equivalence relations on types, there are still some things lacking and some things that are rather strange.

Function extensionality To prove properties about functions, it is often useful to have the principle of function extensionality:

$$\begin{aligned} \text{functionExtensionality} : (A \ B : \text{Type}) \rightarrow (f \ g : A \rightarrow B) \\ \rightarrow ((x : A) \rightarrow f \ x \equiv g \ x) \\ \rightarrow f \equiv g \end{aligned}$$

However, in Martin-Löf's type theory there is no term of that type. Since this theory has the canonicity property, having a propositional equality in the empty context, i.e. $\vdash p : x \equiv y$, we know that p must be canonical: it is definitionally equal to *refl*. In order for it to type check, we then know that x and y must be definitionally equal. Now consider the functions $f = \lambda n \rightarrow n + 0$ and $g = \lambda n \rightarrow 0 + n$, with the usual definition of $+$: $\mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$, we can prove that $(n : \mathbb{N}) \rightarrow f \ n \equiv g \ n$, but not that $f \equiv g$, since that would imply they are definitionally equal, which they are not.

Uniqueness of identity proofs The canonicity property implies that if, in the empty context, we have two identity proofs $p \ q : \text{Id } A \ x \ y$, these proofs are both *refl*, hence they are equal to one another. One would expect that it is possible to prove this inside Martin-Löf's type theory. Using dependent pattern matching, we can easily prove this property in Agda, called uniqueness of identity proofs :

$$\begin{aligned} \text{UIP} : (A : \text{Type}) (x \ y : A) (p \ q : \text{Id } A \ x \ y) \rightarrow \text{Id } (\text{Id } A \ x \ y) \ p \ q \\ \text{UIP } A \ x \ .x \ \text{refl} \ \text{refl} = \text{refl} \end{aligned}$$

Proving this using *J* instead of dependent pattern matching to prove uniqueness of identity proofs has remained an open problem for a long time and has eventually been shown to be impossible [Hofmann and Streicher, 1996] by constructing a model of Martin-Löf's type theory in which there is a type that violates

uniqueness of identity proofs. This tells us that dependent pattern matching is a non-conservative extension over Martin-Löf's type theory¹

As a complement to J , Streicher introduced the induction principle K :

$$\begin{aligned} K &: (A : \text{Type}) (x : A) (P : \text{Id } A \ x \ x \rightarrow \text{Type}) \\ &\rightarrow P \ \text{refl} \\ &\rightarrow (c : \text{Id } A \ x \ x) \\ &\rightarrow P \ c \end{aligned}$$

Using K we can prove UIP , and the other way around. It has been shown that in type theory along with axiom K , we can rewrite definitions written with dependent pattern matching to ones that use the induction principles and axiom K [Goguen et al., 2006].

2.3 Homotopy interpretation

In the introduction (chapter 1), it was mentioned that homotopy type theory concerns itself with the following correspondence:

type theory	homotopy theory
A is a type	A is a space
$x, y : A$	x and y are points in A
$p, q : x \equiv y$	p and q are paths from x to y
$w : p \equiv q$	w is a homotopy between paths p and q

In section 2.1 we noted that homotopies have a ∞ -groupoid structure. It is this structures that leads us to the correspondence between the identity types from Martin-Löf's type theory and homotopy theory. In Hofmann and Streicher [1996], the authors note that types have a groupoid structure. We have a notion of composition of proofs of propositional equality: the term $\text{trans} : \text{Id } A \ x \ y \rightarrow \text{Id } A \ y \ z \rightarrow \text{Id } A \ x \ z$, as such we will use the notation $_ \circ _$ instead of trans . The same goes for $\text{symm} : \text{Id } A \ x \ y \rightarrow \text{Id } A \ y \ x$, which we will denote as $_^{-1}$. We can prove that this gives us a groupoid, i.e. we can prove the following laws hold:

Given $a, b, c, d : A$ and $p : a \equiv b$, $q : b \equiv c$ and $r : c \equiv d$ we have:

- Associativity: $p \circ (q \circ r) \equiv (p \circ q) \circ r$
- Left inverses: $p^{-1} \circ p \equiv \text{refl}$
- Right inverses: $p \circ p^{-1} \equiv \text{refl}$
- Left identity: $\text{refl} \circ p \equiv p$
- Right identity: $p \circ \text{refl} \equiv p$

¹This actually means that all the code we write, should be written using the elimination principles. Agda provides a `WITHOUTK` flag that limits pattern matches to those cases that should be safe. The assumption is that every definition given by pattern matching that passes the `WITHOUTK` check, can be rewritten using the elimination principles. As such, we will sometimes use pattern matching for our definition.

The important thing to note is what kind of equalities we were talking about: associativity, etc. all hold up to propositional equality one level higher. The identity type $Id\ A\ x\ y$ is of course a type and therefore has a groupoid structure of its own. Every type gives rise to a tower of groupoids that can interact with each other. This is exactly the same as the way homotopies form an ∞ -groupoid, hence we have the correspondence between types and spaces as mentioned earlier.

Having such an interpretation of type theory brings us several things. Since every proof we write in type theory corresponds to a proof of a statement in homotopy theory, we can use it to prove theorems of homotopy theory.

It also means that the intuition about homotopy theory can be applied to type theory. As such, we can use it to explain why one cannot prove K using J (subsection 2.3.1), using a couple of illustrations.

2.3.1 Interpreting uniqueness of identity proofs and K

Recall the elimination principle J :

$$\begin{aligned} J : & (A : Type) \\ & \rightarrow (x : A) \\ & \rightarrow (P : (y : A) \rightarrow (p : Id\ A\ x\ y) \rightarrow Type) \\ & \rightarrow (c : P\ x\ x\ refl) \\ & \rightarrow (y : A) \rightarrow (p : Id\ A\ x\ y) \\ & \rightarrow P\ x\ y\ p \end{aligned}$$

Interpreting J in homotopy theory, we see that it tells us that if we want to prove a property about a predicate P on paths, we only have to show that it holds for the constant path $refl$. Homotopically this can be motivated by the fact that P is a predicate on paths with a fixed starting point x and a y that can be chosen freely (see subsection 2.3.1). Any path $p : x \equiv y$ can be contracted along this path to the constant path $refl : x \equiv x$.

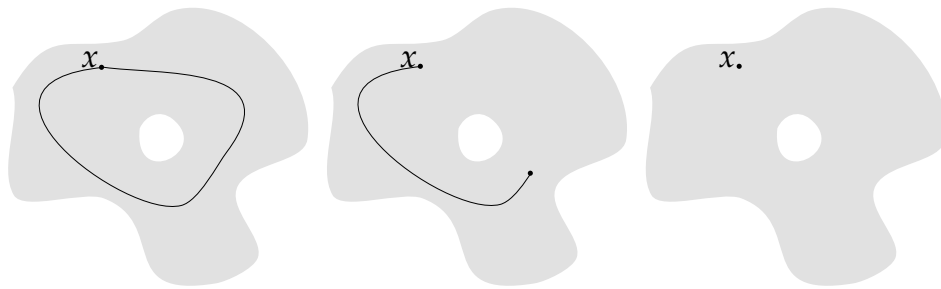


Figure 2.2: With J we have the freedom to move the end point around.

In the case of axiom K , both the beginning and the end point are fixed:

$$\begin{aligned} K : & (A : Type) (x : A) (P : Id\ A\ x\ x \rightarrow Type) \\ & \rightarrow P\ refl \\ & \rightarrow (p : Id\ A\ x\ x) \\ & \rightarrow P\ c \end{aligned}$$

Homotopically this means that we are restricted to loops. If we want to contract a given path $p : x \equiv x$ to $\text{refl} : x \equiv x$, we cannot use the same trick as with J , as the end point is fixed. Contracting any loop to refl does not always work, as can be seen in Figure 2.3.1. If we have a hole in our space, then we can distinguish between loops that go around the hole and those that do not.

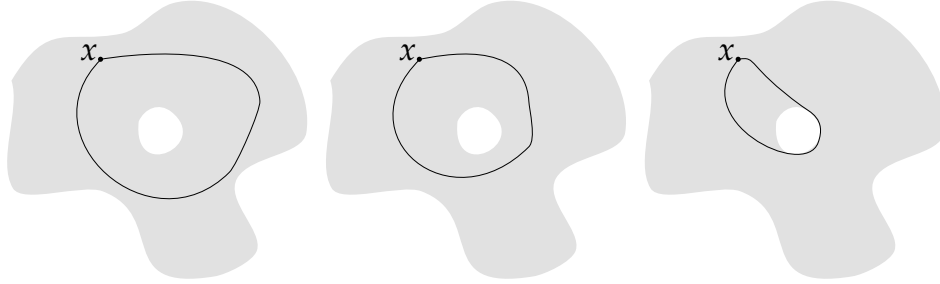


Figure 2.3: With K , we are restricted to loops

2.4 n -types and truncations

The tower of iterated identity types of a type can tell us all sorts of things about the type. For example, we can have a tower in which the identity types in a sense become simpler every iteration, until they reach a fixpoint, in which the identity types are isomorphic to the unit type, \top . In homotopy theory, spaces isomorphic (or rather, homotopic) to the “unit space”, i.e. the space consisting of one point, are called contractible. One way to formulate this in type theory is with the following definition:

$$\begin{aligned} \text{isContractible} &: \text{Type} \rightarrow \text{Type} \\ \text{isContractible } A &= \Sigma A (\lambda \text{center} \rightarrow (x : A) \rightarrow (\text{Id } A \text{ center } x)) \end{aligned}$$

If the structure of the identity types peters out after n iterations, we call such a type an $(n - 2)$ -type, or $(n - 2)$ -truncated²:

$$\begin{aligned} \text{is-truncated} &: \mathbb{N}_{-2} \rightarrow \text{Type} \rightarrow \text{Type} \\ \text{is-truncated } (-2) \ A &= \text{isContractible } A \\ \text{is-truncated } (S \ n) \ A &= (x \ y : A) \rightarrow \text{is-truncated } n \ (\text{Id } A \ x \ y) \end{aligned}$$

These truncation levels have the property that every n -type is also an $(n + 1)$ -type, i.e. is-truncated defines a filtration on the universe of types.

The *contractible* types are the types that are isomorphic to \top in the sense that a contractible type has an inhabitant that is unique up to propositional equality. In section 2.5 we will see examples of contractible types that have more than one canonical element.

²The somewhat strange numbering, starting at -2 comes from homotopy theory, where they first considered groupoids without any higher structure to be 0-truncated and then generalised backwards.

Types of truncation level -1 are called *h*-propositions. (-1) -types are either empty (\perp) or, if they are inhabited, contractible, hence isomorphic to \top . One can easily prove that *h*-propositions satisfy the principle of proof irrelevance:

```
proofIrrelevance : Type → Type
proofIrrelevance A = (x y : A) → Id A x y
```

The converse also holds: if a type satisfies proof irrelevance, it is an *h*-proposition. Showing this is a bit more involved, but it is a nice example of how one can prove things about equalities between equalities.

```
proofIrrelevance⇒is-proposition : (A : Type) → (p : proofIrrelevance A) → is-hProp A
```

We need to show that for every $x y : A$, $x \equiv y$ is contractible: we need to find a proof $c : x \equiv y$ and show that any other proof of $x \equiv y$ is equal to c . An obvious candidate for c is $p \ x \ y$. To show that $c \equiv p \ x \ y$, we use (one-sided) induction on c , fixing the y , so we need to prove that $refl \equiv p \ y \ y$. Instead of doing this directly, we first prove something more general:

```
lemma : (x y : A) (q : x ≡ y) → p x y ≡ q ∘ p y y
```

This can be done by one-sided induction on q , fixing y . The goal then reduces to showing that $p \ y \ y \equiv p \ y \ y$. Using the lemma we can show that $p \ y \ y \equiv p \ y \ y \circ p \ y \ y$. Combining this with $p \ y \ y \circ refl$ and the fact that $\lambda q \rightarrow p \circ q$ is injective for any p , we get that $p \ y \ y \equiv refl$.

The definition of *h*-proposition fits the classical view of propositions and their proofs: we only care about whether or not we have a proof of a proposition and do not distinguish between two proofs of the same proposition.

Another important case are the 0-types, also called *h*-sets, which are perhaps the most familiar to programmers. These are the types of which we have that any two inhabitants x and y are either equal to each other in a unique way, or are not equal, i.e. *h*-sets are precisely those types that satisfy uniqueness of identity proofs. The simplest example of a type that is a *h*-set, but not a *h*-proposition is the type *Bool*:

```
data Bool : Type where
  True  : Bool
  False : Bool
```

In fact, most types one defines in Agda are *h*-sets. One characteristic of *h*-sets is given by Hedberg's theorem, which states that every type that has decidable equality (i.e. $(x y : A) \rightarrow x \equiv y + (x \equiv y \rightarrow \perp)$) also is an *h*-set. The only way to define a type that is not an *h*-set in Agda, is to add extra propositional equalities to the type by adding axioms. This is the subject of section 2.5.

Cite Nicolai Kraus paper

Notation Sometimes we will use the notation $A : Prop$ to indicate that A is a type that is an *h*-proposition. In an actual implementation *Prop* would be defined as $\Sigma (A : Type) (is - truncated (-1) A)$. When we refer to A , we are usually not interested in an inhabitant of the Σ -type, but in the first field of that inhabitant, i.e. the $A : Type$. The same holds for the notation $A : Set$.

2.4.1 Truncations

It may happen that we sometimes construct a type of which the identity types have too much structure, e.g. it is a 2-type but we want it to be a 0-type. In homotopy type theory, we have a way to consider a type as though it were an n -type, for some n we have chosen ourselves, the so called n -truncation of a type. Special cases that are particularly interesting are the (-1) -truncation, i.e. we force something to be a h -proposition, which is particularly useful when we want to do logic, and 0-truncation, i.e. we force something to be a h -set. The idea is that we add enough extra equalities to the type such that the higher structure collapses. This can be done using higher inductive types (section 2.5). The general construction is rather involved and not of much interest for the purposes of this thesis.

2.5 Higher inductive types

We have seen a counterexample of a space in which the interpretation of K and uniqueness of identity proofs fails: a space with a hole in it. The question is then if we can construct such counterexamples in the type theory itself. Since we are asking for a type $A : \text{Type}$ for which there is an inhabitant $x : A$ with a non-canonical term $p : \text{Id } A \ x \ x$, we know that we cannot do this in normal Martin-Löf's type theory as this would violate the canonicity property.

Higher inductive types extend inductive types with the possibility add *path constructors* to the definition of a type: instead of giving constructors for the points of a space, we may also give constructors for paths between points, and paths between paths, and so on. Using higher inductive types we can now describe familiar spaces, such as the circle (see also Figure 2.4):

```
data Circle : Type where
  base : Circle
  loop : base ≡ base
```

Apart from defining how we can construct equalities between inhabitants of the type, we also need to specify the elimination principle. Roughly speaking we need to ensure that all the points get mapped in such a way that all the equalities are respected. In the case of the circle this looks as follows:

```
Circle → rec : { B : Set }
  → (b : B)
  → (p : b ≡ b)
  → Circle → B
```

with computation rule:

```
Circle → rec b p base = b
```

We also need a computation rule for the paths, to witness that the *loop* indeed gets mapped onto the specified path $p : b \equiv b$ by ap :

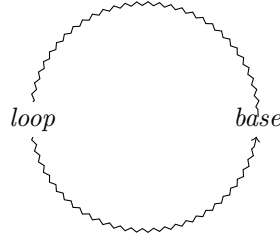


Figure 2.4: The circle as a higher inductive type

$$ap \ (Circle - rec \ b \ p) \ loop = p$$

It might seem a bit silly that we need to provide a path $b \equiv b$, as this type is always inhabited by $refl$. However, we sometimes do want p to be different from $refl$: in order to write the identity function on $Circle$, we also want $loop$ to be preserved by this map.

Apart from a non-dependent elimination principle, we also need a dependent version:

$$\begin{aligned} Circle - ind : \{ B : Circle \rightarrow Set \} \\ \rightarrow (b : B \ base) \\ \rightarrow (p : transport \ B \ loop \ b \equiv b) \\ \rightarrow (x : Circle) \rightarrow B \ x \end{aligned}$$

Using the dependent elimination principle, we can show that this type violates uniqueness of identity proofs, i.e. we can prove that $loop \equiv refl$ does not hold. In fact, the $Id \ Circle \ base \ base$ is isomorphic to the integers \mathbb{Z} , where transitivity maps to addition on integers [Licata and Shulman, 2013]. This might seem a bit strange, because at first glance $Circle$ seems to be a contractible type: we have only one constructor $base$ and an equality $base \equiv base$, so it seems to fit the definition. However, trying to prove $(x : Circle) \rightarrow x \equiv base$ will not work, as the only functions we can define in type theory are *continuous* functions. While it is true in homotopy theory that for every point on the circle, we can find a path to the base point, we cannot do so in a continuous way.

If we add a path constructor connecting two points x and y , we do not only get that specific path, but all the paths that can be constructed from that path using transitivity and symmetry. If we start out with a type with only two constructors x and y , we get a type isomorphic to the booleans (see Figure 2.5), a 0-type. Adding one path constructor $p : x \equiv y$ gives us the interval (see Figure 2.6 and subsection 2.5.2), which is a contractible type (it is a (-2) -type) and hence isomorphic to the unit type \top . If we add yet another path constructor $q : x \equiv y$, we get a type isomorphic to $Circle$, which is a 1-type.

2.5.1 Coherence issues

Equalities at different levels interact with each other: if we add equalities at one level, e.g. paths between points, it may also generate new paths at other levels,

x y

Figure 2.5: Booleans

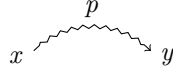


Figure 2.6: Interval

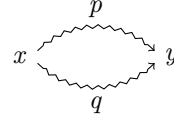


Figure 2.7: Circle

e.g. new homotopies between paths that previously did not exist. One example of this is (-1) -truncation, or propositional truncation, via the following higher inductive type:

```
data (-1)-truncation : (A : Type) : Type where
  inhabitant : A → (-1)-truncation A
  all-paths : (x y : (-1)-truncation A) → x ≡ y
```

We have seen in section 2.4 that this type indeed yields a proposition, as it satisfies proof irrelevance. It collapses all higher structure of the original type $A : Type$.

The converse can also happen: instead of collapsing the structure at higher levels, we might gain new structure at those levels, which sometimes may be undesirable. Suppose we want to consider words generated by some alphabet $A : Set$. This can be done with the following type:

```
data FreeSemigroup : (A : Set) : Type where
  elem : A → FreeSemigroup A
  _·_ : FreeSemigroup A → FreeSemigroup A → FreeSemigroup A
```

Clearly, $FreeSemigroup A$ is a set. Suppose we want $_·_$ to be associative, so we add the following path constructor:

```
assoc : {a b c : FreeSemigroup A} → (a · b) · c ≡ a · (b · c)
```

Adding these equalities breaks the set property: the following diagram (the so called *Mac Lane pentagon*) does not commute:

cite

$$\begin{array}{ccc}
 ((a \cdot b) \cdot c) \cdot d & \xrightarrow{ap (\lambda x \rightarrow x \cdot d) \text{ assoc}} & (a \cdot (b \cdot c)) \cdot d \xrightarrow{\text{assoc}} a \cdot ((b \cdot c) \cdot d) \\
 \downarrow \text{assoc} & & \downarrow ap (\lambda x \rightarrow a \cdot x) \text{ assoc} \\
 (a \cdot b) \cdot (c \cdot d) & \xrightarrow{\text{assoc}} & a \cdot (b \cdot (c \cdot d))
 \end{array}$$

This shows us that the interaction of propositional equalities at the different levels can be quite subtle. For this reason one often truncates a higher inductive type, to be sure that it is coherent enough, e.g. that it is really a h -set.

2.5.2 Interval

Another example of a space from homotopy theory is the interval. At first glance this might seem like a rather uninteresting space to study, as it is homotopy equiv-

alent to the space that consists of one point. The following presentation of the interval as a higher inductive type has some interesting consequences.

The interval $[0, 1]$ can be seen, from a homotopy theory perspective, as a space with two points, 0 and 1, and a path between them. As a higher inductive type, this can be presented as follows:

```
data Interval : Type where
  zero : Interval
  one  : Interval
  segment : zero  $\equiv$  one
```

A map from *Interval* to some type $B : \text{Type}$ must map *zero* and *one* to points in $a, b : B$ such that $a \equiv b$:

```
Interval - rec : { B : Type }
  → (b0 b1 : B)
  → (p : b0  $\equiv$  b1)
  → Interval → B
```

with computation rules:

```
Interval - rec b0 b1 p zero = b0
Interval - rec b0 b1 p one  = b1
ap (Interval - rec b0 b1 p) seg = p
```

Having an interval type means that we have a different way to talk about equalities: any path $p : \text{Id } A \ x \ y$ can be seen as a map $\text{Interval} \rightarrow A$:

```
eqtointerval : { A : Type } { x y : A } → x  $\equiv$  y → Interval → A
eqtointerval { A } { x } { y } p i = Interval - rec { A } x y p i
```

The other way around can also be done:

```
intervaltoeq : { A : Set } → (p : Interval → A) → (p zero)  $\equiv$  (p one)
intervaltoeq p = ap p seg
```

Using this we can now manipulate propositional equalities in such a way that we can prove function extensionality. Suppose two functions $f, g : A \rightarrow B$ and a term $\alpha : (x : A) \rightarrow f \ x \equiv g \ x$. To remove the dependency in the type, we can use *eqtointerval*:

```
 $\lambda a \rightarrow \text{eqtointerval } (\alpha a) : A \rightarrow \text{Interval} \rightarrow A$ 
```

If we flip the arguments of that term, we get a function $\text{Interval} \rightarrow A \rightarrow A$, which then can be turned into the desired $f \equiv g$. The whole term looks as follows:

```
ext : (A B : Type) (f g : A → B) (alpha : (x : A) → f x  $\equiv$  g x) → f  $\equiv$  g
ext A B f g alpha = intervaltoeq (flip ( $\lambda a \rightarrow \text{eqtointerval } (\alpha a)$ ))
```

2.6 Equivalence and univalence

Martin-Löf's type theory satisfies the property that everything you construct in the theory is invariant under isomorphism. Consider for example the definition of a monoid:

$$\begin{aligned}
 \text{Monoid} &: \text{Type} \\
 \text{Monoid} &= \Sigma (\text{carrier} : \text{Set}) . \\
 &\quad \Sigma (\text{unit} : \text{carrier}) . \\
 &\quad \Sigma (\cdot : \text{carrier} \rightarrow \text{carrier} \rightarrow \text{carrier}) . \\
 &\quad \Sigma (\text{assoc} : (x \ y \ z : \text{carrier}) \rightarrow x \cdot (y \cdot z) \equiv (x \ \text{bin op} \ y) \cdot z) . \\
 &\quad \Sigma (\text{unitleft} : (x : \text{carrier}) \rightarrow \text{unit} \cdot x \equiv x) . \\
 &\quad \Sigma (\text{unitright} : (x : \text{carrier}) \rightarrow x \cdot \text{unit} \equiv x) . \top
 \end{aligned}$$

If we have two types $A \ B : \text{Type}$ with an isomorphism $f : A \rightarrow B$ and a proof $ma : \text{Monoid } A$, then it is straightforward to produce a $\text{Monoid } B$ using only $\text{Monoid } A$ and the isomorphism f , by applying f and f^{-1} to the fields of $ma : \text{Monoid } A$. The resulting instance of $\text{Monoid } B$ can then also be shown to be isomorphic to ma . This is similar to the situation with *transport* and *apd*: if we have a proof $p : A \equiv B$, then we can use *transport* to create an inhabitant of $\text{Monoid } B$ using ma and p . We can then prove that the resulting instance of $\text{Monoid } B$ is propositionally equal to ma using *apd*. However, writing *transport* and *apd* function that works with isomorphisms instead of propositional equalities will not work in Martin-Löf's type theory, as we cannot access the information about how the types are constructed, to figure out where the isomorphisms have to be applied.

Univalence gives us an internal account of this principle. It roughly says that isomorphic types are propositionally equal, so all the tools to manipulate propositional equalities now also can be applied to isomorphisms. But before we can formulate the univalence axiom, we need to introduce some new terminology. We can define the notion of a function $f : A \rightarrow B$ being an isomorphism as follows:

$$\begin{aligned}
 \text{isIsomorphism} &: \{A \ B : \text{Type}\} (f : A \rightarrow B) \rightarrow \text{Type} \\
 \text{isIsomorphism } f &= \Sigma (B \rightarrow A) (\lambda g \rightarrow (x : B) \rightarrow f (g \ x) \equiv x \times \\
 &\quad (x : A) \rightarrow g (f \ x) \equiv x)
 \end{aligned}$$

We want the type $\text{isIsomorphism } f$ to be an h -proposition, which it is when A and B are h -sets, but it can fail to be an h -proposition when A and B are n -types with $n > 0$. Instead we introduce the notion of *equivalence*:

$$\begin{aligned}
 \text{isEquivalence} &: \{A \ B : \text{Type}\} (f : A \rightarrow B) \rightarrow \text{Type} \\
 \text{isEquivalence } f &= \Sigma (B \rightarrow A) (\lambda g \rightarrow (x : B) \rightarrow f (g \ x) \equiv x) \\
 &\quad \times \quad \Sigma (B \rightarrow A) (\lambda h \rightarrow (x : a) \rightarrow h (f \ x) \equiv x)
 \end{aligned}$$

This definition does satisfy the property that $\text{isEquivalence } f$ can hold in at most one way (up to propositional equality). We can also show that $\text{isIsomorphism } f \rightarrow \text{isEquivalence } f$ and $\text{isEquivalence } f \rightarrow \text{isIsomorphism } f$, i.e. the two types are *coinhabited*.

Using this definition of what it means to be an equivalence, we can define the following relation on types:

$$\begin{aligned} \simeq & : \{ A B : Type \} \rightarrow Type \\ A \simeq B & = \Sigma (A \rightarrow B) (\lambda f \rightarrow isEquivalence f) \end{aligned}$$

It is easy to show that if two types are propositional equal, then they are also equivalent, by transporting along $\lambda X \rightarrow X$:

$$\equiv \Rightarrow \simeq : (A B : Type) \rightarrow A \equiv B \rightarrow A \simeq B$$

The *univalence axiom* then tells us that equivalences and propositional equalities are equivalent:

$$Univalence : (A B : Type) \rightarrow isEquivalence (\equiv \Rightarrow \simeq A B)$$

One important consequence of this axiom is that we have the following:

$$univalence : (A B : Type) \rightarrow A \simeq B \rightarrow A \equiv B$$

which should satisfy the following computation rule:

$$\begin{aligned} uacomp & : \{ A B : Type \} \\ & \{ f : A \rightarrow B \} \\ & \{ eq : isEquivalence f \} \\ & \{ x : A \} \\ & \rightarrow transport (\lambda X \rightarrow X) (univalence A B) x \equiv f x \end{aligned}$$

Univalence means that we can now generalise the *Monoid* example mentioned, since *transport* and *apd* can now be used for isomorphisms as well.

If we have univalence, the universe of *h*-sets is not a *h*-set, as is exhibited by the isomorphisms $Bool \rightarrow Bool$. There are two different such isomorphisms: *id* and *not*. Using *univalence*, these isomorphisms map to different proofs of $Bool \equiv Bool$. *id* maps to *refl* and *not* to something that is not equal to *refl*. This means that the universe of *h*-sets violates uniqueness of identity proofs. It can be shown to be a 1-type instead. In fact, the universe of *n*-type is not an *n*-type but an $(n + 1)$ -type [Kraus and Sattler, 2013].

2.7 Implementation

Currently, the way to “implement” homotopy type theory, i.e. Martin-Löf’s type theory with univalence and higher inductive types, is to take an existing implementation of Martin-Löf’s type theory such as Agda or Coq and adding univalence and the computation rules for univalence as axioms. This approach is sufficient when we want to do formal mathematics, since in that case we only are interested in type checking our developments. If we want to run the program, terms that make use of univalence then get stuck as soon as it hits an axiom.

The computational interpretation of univalence is one of the biggest open problems of homotopy type theory. Several attempts have been made at a computational interpretation for truncated versions of homotopy type theory: [Licata and Harper \[2012\]](#) show that if we restrict ourselves to a univalent universe of h -sets, we can achieve canonicity. The article however does not present a decidability result for type checking. [Sozeau et al. \[2013\]](#) internalise homotopy type theory in Coq and also restrict themselves to the two-dimensional case, i.e. uniqueness of identity proofs need not hold, but equalities between equalities are unique.

The conjecture is that full canonicity will probably not hold, but only canonicity “up to propositional equality”: it is conjectured [\[Voevodsky, 2011\]](#) that there is a terminating algorithm that takes an expression $t : \mathbb{N}$ and produces a canonical term $t' : \mathbb{N}$ along with a proof that $t \equiv t'$. The proof of equality may use the univalence axiom.

higher inductive types can also be implemented by adding axioms for the extra paths. The elimination principles also can be implemented by adding the computation rules for paths as axioms. One then has to be careful to not do pattern matching on higher inductive types. In Agda one can hide things in such a way that one can export an elimination principle in which the computation rules for the points hold definitionally and the other rules propositionally, while also making direct pattern matching impossible from any other module that imports the module containing the higher inductive type [\[Licata, 2011\]](#). However, one still has to be careful not to use the absurdity pattern, `()`, when dealing with higher inductive types, as that can be used to prove \perp [\[Danielsson, 2012\]](#).

Since dependent pattern matching is not a conservative extension of Martin-Löf’s type theory and in general incompatible with homotopy type theory, we have to use the `--without-K` flag for our Agda code, to ensure that we aren’t pattern matching too liberally. The assumption is that all code written using pattern matching that passes the `--without-K` check can be rewritten using only elimination principles.

A question one might ask is whether we cannot add an extra constructor to the definition of `Id` for univalence. Doing this means that we end up with a different elimination principle: if we want to prove something about propositional equalities, we also need to account for the case when it was proven using univalence. Apart from making it more difficult to prove things about propositional equalities, it is also not clear if the resulting type has the right properties to be called an equality, if we look at its interpretation in some model.

Chapter 3

Applications of homotopy type theory

In chapter 2 we introduced homotopy type theory and two extensions to Martin-Löf’s type theory, univalence and higher inductive types. We have seen how higher inductive types can be used to prove function extensionality and how univalence makes it a lot easier to deal with isomorphic types. This chapter is devoted to other applications of homotopy type theory to programming. In section 3.1 we show how higher inductive types can be used to define quotient types, argue whether we need such a construction and contrast this approach to the setoid approach. We consider some of the difficulties that higher inductive types usually bring with them (so called coherence issues) and show how to write binary operations on quotients as an example of how one uses the elimination principles.

We also consider the application of univalence to views on abstract types (section 3.2), as proposed by [Licata \[2012\]](#). We work out the computations in detail to show how this works out and extend the approach to also work with non-isomorphic views.

3.1 Quotient types

In mathematics, one way to construct new sets is to take the *quotient* of a set X by an equivalence relation R on that particular set. The new set is formed by regarding all elements $x, y \in X$ such that xRy as equal. An example of a quotient set is the set of rationals \mathbb{Q} constructed from the integers as follows: we quotient out $\mathbb{Z} \times \mathbb{Z}$ by the relation $(a, b) \sim (c, d)$ if and only if $ad = bc$.

In programming, such a construction can also be very useful, as it often happens that we have defined a data type that has more structure than we want to expose. This situation typically occurs when we want to encode our data in such a way that certain operations on the data can be implemented more efficiently. An example of this is implementing a dictionary with a binary search tree: there are

multiple binary search trees that represent the same dictionary, i.e. contain the same key-value pairs. If we pass two different trees representing the same dictionary to an operation, we want the operation to yield the same results.

To make the above more precise, suppose we have defined a data type of binary search trees, $BST : Type$, along with a relation $\sim : BST \rightarrow BST \rightarrow hProp$ such that $x \sim y \equiv \top$ if and only if x and y are comprised of the same key-value pairs, and $x \sim y \equiv \perp$ otherwise. Suppose we have an insertion operation $insert : KeyValuePair \rightarrow BST \rightarrow BST$ and a lookup function $lookup : Key \rightarrow BST \rightarrow Maybe Value$. We can formulate the properties that should hold:

- $(a : KeyValPair) (x y : BST) \rightarrow x \sim y \rightarrow insert\ a\ x \sim insert\ a\ y$
- $(a : Key) (x y : BST) \rightarrow x \sim y \rightarrow lookup\ a\ x \equiv lookup\ a\ y$

Note that for insertion, returning the same results means that we want them to represent the same dictionary: it is perfectly allowed to return differently balanced binary search trees. For *lookup*, we want the results to be propositionally equal, as we do not have any other relation available that holds on the result type, *Maybe Value*.

A type that comes equipped with an equivalence relation, such as BST along with \sim , is called a *setoid*. Its disadvantages are that we have to formulate and check the properties ourselves: there is no guarantee that a function out of a setoid respects the relation from the setoid. As can be seen in the binary search tree example, we have to be careful to use the right relation (propositional equality or the setoid's equivalence relation) when we want to talk about two inhabitants being the same. Homotopy type theory provides us with the machinery, namely higher inductive types, to enrich the propositional equality of a type, so we can actually construct a new type in which propositional equality and the provided equivalence relation coincide.

3.1.1 Do we need quotients?

Before we look at the quotient type construction with higher inductive types, we will determine whether we actually need such a thing. In the case of the dictionary example, we might consider making the BST data type more precise such that the only inhabitants are trees that are balanced in a certain way, so we do have a unique representation for every dictionary.

The question then is whether such a construction always exists: can we define a type that is in some sense equal to the quotient? To be able to answer this question, we need to define what it means to be a quotient and what notion of equality we want.

[Altenkirch et al.](#) define a quotient, given a setoid (A, \sim) as a type $Q : Type$ with the following:

- a projection function $[_] : A \rightarrow Q$
- a function $sound : (x y : A) \rightarrow x \sim y \rightarrow [x] \equiv [y]$
- an elimination principle:

$$\begin{aligned}
Q\text{-elim} : & (B : Q \rightarrow \text{Type}) \\
& (f : (x : A) \rightarrow B [x]) \\
& ((x \ y : A) (p : x \sim y) \rightarrow (\text{transport } (\text{sound } x \ y \ p) (f \ x)) \equiv f \ y) \\
& (q : Q) \rightarrow B \ q
\end{aligned}$$

A quotient is called definable if we have a quotient Q along with the following:

- $emb : Q \rightarrow A$
- $complete : (a : A) \rightarrow emb [a] \sim a$
- $stable : (q : Q) \rightarrow [emb \ q] \equiv q$

We can view these requirements as having a proof of $[_]$ being an isomorphism, with respect to the relation \sim on A instead of propositional equality.

The result of [Altenkirch et al.](#) is that there exist quotients that are not definable with one example being the real numbers constructed using the usual Cauchy sequence method. Adding quotients as higher inductive types to our type theory, does not make the real numbers definable. Adding quotients is still useful in that we only have to work with propositional equality, as opposed to the confusion as to what relation one should use that arises from the use of setoids.

3.1.2 Quotients as a higher inductive type

Using higher inductive types, we can define the quotient of a type by a relation as follows:

$$\begin{aligned}
\text{data } Quotient \ (A : \text{Type}) \ (\sim : A \rightarrow A \rightarrow hProp) : \text{Type} \text{ where} \\
\quad [-] : A \rightarrow Quotient \ A \ \sim \\
\quad \text{sound} : (x \ y : A) \rightarrow x \sim y \rightarrow [x] \equiv [y]
\end{aligned}$$

To write a function $Quotient \ A \ \sim \rightarrow B$ for some $B : \text{Type}$, we need to specify what this function should do with values $[x]$ with $x : A$. This needs to be done in such a way that the paths added by sound are preserved. Hence the recursion principle lifts a function $f : A \rightarrow B$ to $\tilde{f} : Quotient \ A \ \sim \rightarrow B$ given a proof that it preserves the added paths:

$$\begin{aligned}
Quotient\text{-}rec : & (A : \text{Type}) (\sim : A \rightarrow A \rightarrow hProp) \\
& (B : \text{Type}) \\
& (f : A \rightarrow B) \\
& ((x \ y : A) \rightarrow x \sim y \rightarrow f \ x \equiv f \ y) \\
& Quotient \ A \ \sim \rightarrow B
\end{aligned}$$

If we generalise this to the dependent case, we get something that fits perfectly in the requirement of a type being a quotient given earlier:

$$\begin{aligned}
Quotient\text{-}ind : & (A : \text{Type}) (\sim : A \rightarrow A \rightarrow hProp) \\
& (B : Quotient \ A \ \sim \rightarrow \text{Type}) \\
& (f : (x : A) \rightarrow B [x])
\end{aligned}$$

$$\begin{aligned} & ((x\ y : A) (p : x \sim y) \rightarrow (\text{transport } (\text{sound } x\ y\ p) (f\ x)) \equiv f\ y) \\ & (q : \text{Quotient } A\ \sim) \rightarrow B\ q \end{aligned}$$

Note that we do not require a proof of \sim being an equivalence relation. Instead, the quotient should be read as identifying inhabitants by the smallest equivalence relation generated by \sim .

3.1.3 Coherence issues

One thing we glossed over is the question whether $\text{Quotient } A\ \sim$ is actually a h -set. This need not be the case, as is exhibited by the case where A is taken to be \top and \sim is the trivial relation. The resulting quotient is equivalent to the circle, which is not a h -set: the loop $\text{sound } tt\ tt\ tt : [tt] \equiv [tt]$ is not equal to $\text{refl} : [tt] \equiv [tt]$.

In order to get a h -set, we therefore need to take the 0-truncation of the quotient, which can be done with the following higher inductive type:

$$\begin{aligned} & \text{0-truncation } (A : \text{Type}) : \text{Type} \text{ where} \\ & \text{inhabitant} : A \rightarrow \text{0-truncation } A \\ & \text{uip} : \{x\ y : \text{0-truncation } A\} \rightarrow (p\ q : x \equiv y) \rightarrow p \equiv q \end{aligned}$$

The elimination principle tells us that any function $A \rightarrow B$, with $A\ B : \text{Type}$ can be lifted to $\text{0-truncation } A \rightarrow B$ if it respects the additional paths of $\text{0-truncation } A$. If B happens to be a h -set, then these conditions are automatically satisfied. In the dependent case, we have to supply a family of types $B : \text{0-truncation } A \rightarrow \text{Type}$ and a function $f : (x : A) \rightarrow B\ (\text{inhabitant } x)$ such that, again, the additional paths of $\text{0-truncation } A$ are respected. If we have that for every $x : \text{0-truncation } A$, $B\ x$ is a h -set, then we are done. The precise formulation of the elimination principles, both dependent and non-dependent, are rather technical and involved and not of interest for our purposes. In the examples we consider, we eliminate into h -sets, so we do not need to explicitly check the additional conditions.

If the relation \sim happens to be an equivalence relation, using the truncated quotient also gives us that we have (using univalence) $a \sim b \equiv ([a] \equiv [b])$, for every $a\ b : A$.

3.1.4 Binary operations on quotients

We have seen how to lift a function $f : A \rightarrow B$ to $\tilde{f} : \text{Quotient } A\ \sim \rightarrow B$ given a proof of $(x\ y : A) \rightarrow x \sim y \rightarrow f\ x \equiv f\ y$, using Quotient-rec . Suppose we want to write a binary operation on quotients, then we want to have a way to lift a function $f : A \rightarrow A \rightarrow B$ satisfying $(x\ y\ x'\ y' : A) \rightarrow x \sim x' \rightarrow y \sim y' \rightarrow f\ x\ y \equiv f\ x'\ y'$ to $\tilde{f} : \text{Quotient } A\ \sim \rightarrow \text{Quotient } A\ \sim \rightarrow B$.

Let us fix A, \sim and B , so that we do not have to pass them around explicitly. Our goal is to write a term of the following type:

$$\begin{aligned}
& \text{Quotient-rec-2} : (f : A \rightarrow A \rightarrow B) \\
& \quad (\text{resp} : (x \ y \ x' \ y' : A) \rightarrow x \sim x' \rightarrow y \sim y' \rightarrow f \ x \ y \equiv f \ x' \ y') \\
& \quad \text{Quotient } A \ _\sim _ \rightarrow \text{Quotient } A \ _\sim _ \rightarrow B
\end{aligned}$$

We will first use *Quotient-rec* to lift the left argument, i.e. we want to produce a function of type $\text{Quotient } A \ _\sim _ \rightarrow A \rightarrow B$ and then use *Quotient-rec* on this function to achieve our goal. So let us try writing the function that lifts the left argument:

$$\begin{aligned}
& \text{lift-left} : (f : A \rightarrow A \rightarrow B) \\
& \quad (\text{resp} : (x \ y \ x' \ y' : A) \rightarrow x \sim x' \rightarrow y \sim y' \rightarrow f \ x \ y \equiv f \ x' \ y') \\
& \quad \text{Quotient } A \ _\sim _ \rightarrow A \rightarrow B \\
& \text{lift-left } f \ \text{resp } q = \text{Quotient-rec } f \ \text{goal}_0 \ q
\end{aligned}$$

where $\text{goal}_0 : (x \ x' : A) \rightarrow x \sim x' \rightarrow f \ x \equiv f \ x'$. Since we have quotient types, we also have function extensionality¹, hence we can solve this by proving $(x \ x' \ y : A) \rightarrow x \sim x' \rightarrow f \ x \ y \equiv f \ x' \ y$. However, to be able to use *resp*, we also need a proof of $y \sim y$, so if we assume that \sim is an equivalence relation, we can solve this goal.

We can now fill in *lift-left* in the definition of *Quotient-rec-2*:

$$\text{Quotient-rec-2 } f \ \text{resp } q \ q' = \text{Quotient-rec } (\text{lift-left } f \ \text{resp } q) \ \text{goal}_1 \ q'$$

where $\text{goal}_1 : (y \ y' : A) \rightarrow y \sim y' \rightarrow \text{lift-left } f \ \text{resp } q \ y \equiv \text{lift-left } f \ \text{resp } q \ y'$, which can be proven using *Quotient-ind*. We then only have to consider the case where q is of the form $[a]$ for some $a : A$. In that case, $\text{lift-left } f \ \text{resp } q \ y$ reduces to $f \ a \ y$ and $\text{lift-left } f \ \text{resp } q \ y'$ to $f \ a \ y'$. Since we have $y \sim y'$, we again need \sim to be reflexive to get $a \sim a$ so we can use *resp*. We now have the following:

$$\begin{aligned}
& \text{goal}_1 : (y \ y' : A) \rightarrow y \sim y' \rightarrow \text{lift-left } f \ \text{resp } q \ y \equiv \text{lift-left } f \ \text{resp } q \ y' \\
& \text{goal}_1 = \lambda y \ y' \ r \rightarrow \\
& \quad \text{Quotient-ind } (\lambda w \rightarrow \text{lift-left } f \ \text{resp } w \ y \equiv \text{lift-left } f \ \text{resp } w \ y') \\
& \quad (\lambda a \rightarrow \text{resp } a \ y \ a \ y' (\sim \text{-refl } a) \ r) \\
& \quad \text{goal}_2 \\
& \quad q
\end{aligned}$$

Of course, we have still to prove that this respects the quotient structure on q :

$$\begin{aligned}
& \text{goal}_2 : (p : x \sim x') \\
& \quad \text{transport } (\text{sound } x \ x' \ p) (\text{resp } x \ y \ x \ y' (\sim \text{-refl } x) \ r) \equiv \\
& \quad \text{resp } x' \ y \ x' \ y' (\sim \text{-refl } x') \ r
\end{aligned}$$

Note that this equality is of type $\text{Id } (\text{Id } B \ (f \ x \ y) \ (f \ x \ y'))$, which means that if B happens to be a h -set, we can appeal to uniqueness of identity proofs and we are done.

¹We can quotient *Bool* by the trivial relation. Using this, we can perform essentially the same proof of function extensionality as the one that uses the interval type.

It is interesting to see that even though we do not need \sim to be an equivalence relation for the definition of quotient to work, we do find ourselves in need of properties such as reflexivity for \sim . This stems from the fact that the relation we take the quotient by, is the smallest *equivalence* relation generated by \sim .

3.2 Views on abstract types

Consider the dictionary example of the previous section. Most languages provide such a structure as an *abstract type*, e.g. in the Haskell Platform, a dictionary structure is provided by the `Data.Map` module. To the users importing this module, the type `Map` is opaque: its constructors are hidden. The user may only use the operations such as `insert` and `lookup`. The advantage of this approach is that we can easily interchange an obvious but slow implementation (e.g. implementing a dictionary as a list of tuples) with a more efficient but more complex solution (e.g. using binary search trees instead of lists), without having to change a single line of code in the modules using the abstract type.

In dependently typed programming, such an approach often means that we have hidden too much: as soon as we try to prove properties about our program that uses some abstract type, we find ourselves having to add properties to the abstract type specification, or even worse: we end up exporting everything so we can use induction on the concrete type used in the actual implementation.

A solution to this problem is to supply the abstract type along with a concrete implementation of the abstract type, called a *view*. This approach was introduced by Wadler [1987] as a way to do pattern matching on abstract types.

3.2.1 Specifying views

An implementation of an abstract type is a type along with a collection of operations on that type. An abstract type can then be described in type theory as a nested Σ -type [Mitchell and Plotkin, 1988], e.g. a sequence abstract type can be described as follows:

$$\begin{aligned} \text{Sequence} = \Sigma \quad & (seq : Set \rightarrow Set) \\ & \Sigma (empty : (A : Set) \rightarrow (seq A)) \\ & \Sigma (single : (A : Set) \rightarrow A \rightarrow seq A) \\ & \Sigma (append : (A : Set) \rightarrow seq A \rightarrow seq A \rightarrow seq A) . \\ & (map : (A B : Set) \rightarrow (A \rightarrow B) \rightarrow seq A \rightarrow seq B) \end{aligned}$$

An implementation of such an abstract type then is just an inhabitant of this nested Σ -type.

If we want to do more than just use the operations and prove properties about our programs that make use of abstract types, we often find that we do not have enough information in the abstract type specification available to prove the property at hand. One way to address this problem is to add properties to the specification, but it might not at all be clear a priori what properties are interesting and expressive enough to add to the specification.

Another solution, proposed by Licata [2012], is to use views: along with nested Σ -type, we also provide a concrete implementation, i.e. an inhabitant of said Σ -type, called a *view* on the abstract type. The idea is that the concrete view can be used to prove theorems about the abstract type. However, for this to work, we need to make sure that any implementation of the abstract type is also in some sense compatible with the view: the types of both implementations need to be isomorphic and the operations need to respect the isomorphism. To illustrate this, consider we have two sequence implementations:

$$\begin{aligned} \text{ListImpl} &: \text{Sequence} \\ \text{ListImpl} &\triangleq (\text{List}, ([], (\lambda x \rightarrow [x], (-++-, \text{map})))) \\ \text{OtherImpl} &: \text{Sequence} \\ \text{OtherImpl} &= (\text{Other}, (\text{otherEmpty}, (\text{otherSingle}, (\text{otherAppend}, \text{otherMap})))) \end{aligned}$$

We want *List* and *Other* to be “isomorphic”², i.e. we need to write the following terms:

- $\text{to} : (A : \text{Type}) \rightarrow \text{Other } A \rightarrow \text{List } A$
- $\text{from} : (A : \text{Type}) \rightarrow \text{List } A \rightarrow \text{Other } A$
- $\text{fromIsRightInverse} : (A : \text{Type}) (xs : \text{List } A) \rightarrow \text{to} (\text{from } xs) \equiv xs$
- $\text{fromIsLeftInverse} : (A : \text{Type}) (xs : \text{Other } A) \rightarrow \text{from} (\text{to } xs) \equiv xs$

We also want the operations on *Other* to behave in the same way as the operations on *Lists*, i.e. the following properties should be satisfied:

- $\text{to } \text{otherEmpty} \equiv []$
- $(x : A) \rightarrow \text{to} (\text{otherSingle } x) \equiv [x]$
- $(xs \text{ } ys : \text{Other } A) \rightarrow \text{to} (\text{otherAppend } xs \text{ } ys) \equiv \text{to } xs ++ \text{to } ys$
- $(f : A \rightarrow B) (xs : \text{Other } A) \rightarrow \text{to} (\text{otherMap } f \text{ } xs) \equiv \text{map } f (\text{to } xs)$

These properties can be added to the original *Sequence* type. However, it is rather tedious having to formulate these properties for every operation of the abstract type. Since we have specified the abstract type as a Σ -type, we can use propositional equality between these to guide us to the desired properties. We know that in order to prove that $a \equiv b : \Sigma A B$ are propositionally equal, we need to show its fields are propositionally equal as well:

$$\begin{aligned} \Sigma{-} \equiv & : \{ A : \text{Type} \} \{ B : A \rightarrow \text{Type} \} \\ & \{ s \text{ } s' : \Sigma A B \} \\ & (p : \text{fst } s \equiv \text{fst } s') \\ & (q : \text{transport } B \text{ } p (\text{snd } s) \equiv \text{snd } s') \\ \rightarrow & s \equiv s' \end{aligned}$$

If we want to prove that $\text{ListImpl} \equiv \text{OtherImpl}$, then using $\Sigma{-} \equiv$, we first need to show that $\text{List} \equiv \text{Other}$. This can be done by showing that for every $(A : \text{Type})$, we have an isomorphism $\text{to} : \text{Other } A \rightarrow \text{List } A$. Using the univalence axiom and function extensionality, we can then prove our goal, $\text{List} \equiv \text{Other}$. For the second

²*List* and *Other* cannot be isomorphic, as they are not types but type *constructors*.

part of the outermost Σ -type, we need to transport the *snd* of *ListImpl* along the proof of $List \equiv Other$ we just gave and prove it to be propositionally equal to the *snd* of *OtherImpl*. Rather than deal with the fully general *Sequence* where will show how the transporting looks like for the case when we fix the type parameter. This is done so we do not have to deal with function extensionality and only have to use univalence directly once. We consider the following definitions where we fix the type parameter $A : Type$:

$$\begin{aligned} Sequence_A = & \Sigma (seq_A : Set) . \\ & \Sigma (empty_A : seq_A) . \\ & \Sigma (single_A : A \rightarrow seq_A) . \\ & \Sigma (append_A : seq_A \rightarrow seq_A \rightarrow seq_A) . \\ & (map_A : (A \rightarrow A) \rightarrow seq_A \rightarrow seq_A) \end{aligned}$$

with *ListImpl_A* and *OtherImpl_A* defined straightforwardly from the previous definitions. To show that *ListImpl_A* and *OtherImpl_A*, we need to show using univalence that $List\ A \equiv Other\ A$, so the beginning of the proof looks like this:

$$\begin{aligned} spec : & \\ & (from : List\ A \rightarrow Other\ A) \\ & (to : Other\ A \rightarrow List\ A) \\ & \rightarrow Iso\ (List\ A)\ (Other\ A)\ from\ to \\ & \rightarrow ListImpl_A \equiv OtherImpl_A \\ spec\ from\ to\ iso = & \Sigma- \equiv (univalence\ (List\ A)\ (JoinList\ A)\ iso) \\ & (\Sigma- \equiv goal_0 \\ & (\Sigma- \equiv goal_1 \\ & (\Sigma- \equiv goal_2 \\ & goal_3))) \end{aligned}$$

The first goal, *goal₀*, has type $fst\ (transport\ (univalence\ (List\ A)\ (JoinList\ A)\ iso)\ ([], (\lambda x \rightarrow [x], (-+-, map)))) \equiv otherEmpty$. The left hand side of the equation is stuck, as we made use of the univalence axiom. However, we can prove that the first field of *transport* applied to the dependent pair, is *transport* applied to the first field of the dependent pair:

$$\begin{aligned} \Sigma\text{-transport} : & \\ & \{ Ctx : Type \} \\ & \{ A : Ctx \rightarrow Type \} \{ B : (ctx : Ctx) \rightarrow A\ ctx \rightarrow Type \} \\ & \{ ctx\ ctx' : Ctx \} \\ & \{ x : A\ ctx \} \{ y : B\ ctx\ x \} \\ & (pf : ctx \equiv ctx') \rightarrow \\ & fst\ (transport\ (\lambda c \rightarrow \Sigma\ (A\ c)\ (B\ c))\ pf\ (x, y)) \equiv transport\ (\lambda c \rightarrow A\ c)\ pf\ x \end{aligned}$$

If we apply this to *goal₀*, we now need to show that $transport\ (\lambda c \rightarrow c)\ (univalence\ (List\ A)\ (JoinList\ A)\ iso)\ [] \equiv otherEmpty$, which we can further reduce using the “computation” rule for univalence:

$$\begin{aligned} univalence\text{-comp} : & \\ & \{ A\ B : Type \} \\ & \{ from : A \rightarrow B \} \end{aligned}$$

$$\begin{aligned}
& \{ to : B \rightarrow A \} \\
& \{ iso : Iso\ A\ B\ from\ to \} \\
& \{ x : A \} \\
& \rightarrow transport\ (\lambda X \rightarrow X)\ (univalence\ A\ B\ iso)\ x \equiv from\ x
\end{aligned}$$

We have reduced $goal_0$ to the proof obligation $from\ [] \equiv otherEmpty$. We can apply the same steps to the other goals and recover the properties we formulated earlier.

With the current “implementation” of homotopy type theory done by adding things such as univalence as axioms, we have to do all this rewriting by hand, but it is not unthinkable that a lot of this can be automated.

3.2.2 Reasoning with views

If we want to prove a property about our abstract type, we can now only have to prove that it holds for the concrete view. The resulting proof can then be used to show that it also holds for any other implementation of the abstract type.

As an example of this, we will show that the *empty* operation of our sequence type is the (left) unit of *append*. The case for lists is easy, assuming that $++$ only does induction on its left argument:

$$\begin{aligned}
left\text{-}unit\text{-}append & : (xs : List\ A) \rightarrow [] ++ xs \equiv xs \\
left\text{-}unit\text{-}append\ xs & = refl
\end{aligned}$$

The general case of this statement is:

$$(xs : Other\ A) \rightarrow otherAppend\ otherEmpty\ xs \equiv xs$$

which can be established by the following equational reasoning:

$$\begin{aligned}
& xs \\
& \equiv \{ isomorphism \} \\
& \quad from\ (to\ xs) \\
& \equiv \{ []\ \text{is left unit of } ++ \} \\
& \quad from\ ([] ++ to\ xs) \\
& \equiv \{ specification\ of\ otherImpl \} \\
& \quad from\ (to\ otherEmpty ++ to\ xs) \\
& \equiv \{ specification\ of\ otherAppend \} \\
& \quad from\ (to\ (otherAppend\ otherEmpty\ xs)) \\
& \equiv \{ isomorphism \} \\
& \quad otherAppend\ otherEmpty\ xs
\end{aligned}$$

3.2.3 Non-isomorphic views

An implementation of an abstract type sometimes does not turn out to be isomorphic to the concrete view. An example of this is an implementation of sequences via join lists:

```

data JoinList (A : Type) : Type where
  nil  : JoinList A
  unit : A → JoinList A
  join : JoinList A → JoinList A → JoinList A

```

We have a function $to : JoinList A \rightarrow List A$ that maps nil to nil , $unit\ a$ to $[a]$ and interprets $join$ as concatenation of lists. The other way around, $from : List A \rightarrow JoinList A$ can be constructed by mapping every element a of the input list to $unit\ a$ and then using $join$ to concatenate the resulting list of $JoinLists$.

While we do have that $(ls : List A) \rightarrow to\ (from\ ls) \equiv ls$, it is not the case that $(js : JoinList A) \rightarrow from\ (to\ js) \equiv js$, as to is not injective: to and $from$ do not form an isomorphism: $JoinList$ has a finer structure than $List$. If only the first equality holds, but the second does not, to is called a *retraction* with $from$ as its *section*. It still makes sense to use $JoinList$ as an implementation of sequences. The properties that the operations on $JoinLists$ should respect, do not make use of the fact that $from$ and to are isomorphisms; they can still be used for non-isomorphic views.

Since we are only interested in using the $JoinList$ as a sequence and do not care how the inhabitants are balanced, we can take the quotient by the following relation:

```

 $\sim : JoinList A \rightarrow JoinList A \rightarrow Type$ 
 $x \sim y = to\ x \equiv to\ y$ 

```

The type $Quotient\ (JoinList\ A)\ \sim$ is then isomorphic to $List\ A$. This result can be generalised to arbitrary section-retraction pairs between h -sets A and B : given $r : A \rightarrow B$ and $s : B \rightarrow A$ such that $(a : A) \rightarrow s\ (r\ a) \equiv a$, then B is isomorphic to A/\sim where $x \sim y$ is defined as $r\ x \equiv r\ y$. We have a function $A \rightarrow A/\sim$, namely the constructor box and can write a function $A/\sim \rightarrow A$. If we use $Quotient - rec$ for this, we need to supply a function $f : A \rightarrow A$ such that if $r\ x \equiv r\ y$, then also $f\ x \equiv f\ y$. Choosing f to be $\lambda x \rightarrow s\ (r\ x)$ works. The identity function need not work: if it did, r would be injective and would be an isomorphism. Let us name the functions between A and A/\sim $to-A/\sim$ and $from-A/\sim$. Composing these functions with r or s , we get functions between A/\sim and B that give us the desired isomorphism. Proving that this is an isomorphism mostly involves applying the proof that $r\ (s\ x) \equiv x$ in various ways. We also have to invoke the uniqueness of identity proofs property that A/\sim admits (A is a set and \sim is an equivalence relation) for the induction step on A/\sim . The fact that $to-A/\sim$ is a retraction with $from-A/\sim$ as its section can be proved using the same techniques.

To lift the operations on A to operations on A/\sim we simply apply $to-A/\sim$ and $from-A/\sim$ in the right places. Showing that these lifted operations satisfy the conditions that follow from the specification then boils down to conditions that only refer to the operations on A in relation to those on B , as we will demonstrate with the $JoinList$ example. Let us define $JoinList\ A/\sim$ as $JLAquote$ with $x \sim y$ defined as $to\ x \equiv to\ y$. We have the following functions:

- $to : JoinList\ A \rightarrow List\ A$
- $from : List\ A \rightarrow JoinList\ A$

- $\widetilde{to} : JoinList\ A \rightarrow JoinList\ A/\sim$
- $\widetilde{from} : JoinList\ A/\sim \rightarrow JoinList\ A$

The isomorphism between $JoinList\ A/\sim$ and $List\ A$ is witnessed by $to \circ \widetilde{from} : JoinList\ A/\sim \rightarrow List\ A$ and $to \circ \widetilde{from} : List\ A \rightarrow JoinList\ A$. The empty of $JoinList\ A/\sim$ is $\widetilde{to}\ nil$, which means that we need to establish $to\ (\widetilde{from}\ (\widetilde{to}\ nil)) \equiv []$. We can reduce this goal to $to\ nil \equiv []$ via equational reasoning:

$$\begin{aligned}
& to\ (\widetilde{from}\ (\widetilde{to}\ nil)) \\
& \equiv \{ \text{definition } \widetilde{to} \} \\
& to\ (\widetilde{from}\ (box\ nil)) \\
& \equiv \{ \beta\ \text{reduction} \} \\
& to\ (from\ (to\ nil)) \\
& \equiv \{ to / from\ \text{is a retraction / section} \} \\
& to\ nil
\end{aligned}$$

In general we have that $\widetilde{from}\ (\widetilde{to}\ x) \equiv from\ (to\ x)$ holds for any $x : JoinList\ A$. Deriving the property for *single* goes analogously to the derivation above. The rule for *append* is more interesting as we there also need \widetilde{from} in other positions:

$$\begin{aligned}
& to\ (\widetilde{from}\ (\widetilde{to}\ (join\ (\widetilde{from}\ xs)\ (\widetilde{from}\ ys)))) \\
& \equiv \{ \beta\ \text{reduction} \} \\
& to\ (from\ (to\ (join\ (\widetilde{from}\ xs)\ (\widetilde{from}\ ys)))) \\
& \equiv \{ to / from\ \text{is a retraction / section} \} \\
& to\ (join\ (\widetilde{from}\ xs)\ (\widetilde{from}\ ys))
\end{aligned}$$

We end up with having to prove that $(xs\ ys : JoinList\ A/\sim) \rightarrow to\ (join\ (\widetilde{from}\ xs)\ (\widetilde{from}\ ys)) \equiv to\ (\widetilde{from}\ xs)$ which follows from $(xs\ ys : JoinList\ A) \rightarrow to\ (join\ xs\ ys) \equiv to\ xs \uplus to\ ys$.

The above derivation shows us that we might arrive at equations that are a bit less general than the equations we get from if we were to pretend our retraction-section pair is actually an isomorphism.

Non-isomorphic views via definable quotients

It so happens that the quotient A/\sim is definable. We can use the type $\Sigma\ (x : A) . s\ (r\ x) \equiv x$, i.e. restrict A to those inhabitants for which s and r are isomorphisms. The function $box : A \rightarrow \Sigma\ (x : A) . s\ (r\ x) \equiv x$ is then defined by $\lambda x \rightarrow (s\ (r\ x))$, *ap* $s\ (is-retract\ (r\ x))$, where $is-retract : (x : B) \rightarrow r\ (s\ x) \equiv x$.

Notice that for the quotient type we have the $\lambda x \rightarrow s\ (r\ x)$ in the “deconstructor” (i.e. in the function $\widetilde{from} : JoinList\ A/\sim \rightarrow JoinList\ A$) and here we have it in the constructor (i.e. the function box). This stems from the fact that the soundness of quotient types is enforced by the way they are eliminated. It is only there that we have the obligation to show that we respect the relation on the type. With the Σ -type it is more correctness by construction.

From a computational perspective, the first approach with the quotient types is more desirable, as the values of the type do not carry around any correctness proof.

3.3 Conclusion

Higher inductive types allow us to straightforwardly define quotient types. This definition works better than the setoid method in that we no longer have to be careful whether we use the custom equivalence relation or propositional equality: we only have to consider propositional equality. However, as is common with higher inductive types, we have to take the 0-truncation in the definition of a quotient type. This makes the elimination principle more complex to work with, but since virtually any of the types we work with in programming are h -sets, we usually automatically satisfy the extra conditions that the 0-truncation adds to the elimination principle.

Univalence gives us a very clean way to define specifications of abstract types using concrete views. Working with this specification, e.g. trying to prove that a given implementation satisfies the specification, involves a lot of manual fiddling with the computation rules of Σ -types and univalence. Having a computational interpretation of univalence would obviously be of great importance for this method to be useful.

Using quotient types, we can also define a view on an abstract type that is not isomorphic to the concrete type of the reference implementation, but only instead we have a retraction-section pair between the two types. Any retraction-section pair can be turned into an isomorphism, by quotienting out by the retraction. Such a quotient happens to be definable, which means that we do not need the quotient type construction using higher inductive types to do this. However, the higher inductive type construction does yield a definition that is more amenable to the optimisations that will be discussed in ??, as the proofs that the quotient structure is respected only occur in the calls to the elimination principle, instead of occurring in all the terms of type, which is the case with the definable quotient implementation.

Chapter 4

Erasing propositions

When writing certified programs in a dependently typed setting, we can conceptually distinguish between the *program* parts and the *proof* (of correctness) parts. These are sometimes also referred to as the informative and logical parts, respectively. In practice, these two seemingly separate concerns are often intertwined. Consider for example the sorting of lists of naturals: given some predicate $isSorted : List\ \mathbb{N} \rightarrow List\ \mathbb{N} \rightarrow Type$ that tells us whether the second list is a sorted permutation of the first one, we can write a term of the following type:

$$sort : (xs : List\ \mathbb{N}) \rightarrow \Sigma\ (ys : List\ \mathbb{N})\ (isSorted\ xs\ ys)$$

To implement such a function, we need to provide for every list a sorted list along with a proof that this is indeed a sorted version of the input list. At run-time the type checking has been done, hence the proof of correctness has already been verified: we want to *erase* the logical parts.

Types such as $isSorted\ xs\ ys$ are purely logical: we care more about the presence of an inhabitant than what kind of inhabitant we exactly have at our disposal. In section 4.1 we will give more examples of such types, called *propositions*, and how they can occur in various places in certified programs. In sections 4.2 and 4.3 we review the methods Coq and Agda provide us to annotate parts of our program as being propositions. Section 4.4 reviews the concept of *collapsible families* and how we can automatically detect whether a type is a proposition, instead of annotating them ourselves. In section 4.5 we internalise the concept of collapsible families and try to do the same with the optimisation in section 4.6. The internalised version of collapsibility looks like an indexed version of the concept of *h-propositions*. In section 4.7 we investigate if we can use this to devise an optimisation akin to the optimisation based on collapsibility.

4.1 Propositions

In the *sort* example, the logical part $isSorted\ xs\ ys$ occurs in the result as part of a Σ -type. This means we can separate the proof of correctness from the sorting

itself, i.e. we can write a function $sort' : List\ \mathbb{N} \rightarrow List\ \mathbb{N}$ and a proof of the following:

$$sortCorrect : (xs : List\ \mathbb{N}) \rightarrow isSorted\ xs\ (sort'\ xs)$$

The logical part here asserts properties of the *result* of the computation. If we instead have assertions on our *input*, we cannot decouple this from the rest of the function as easily as, if it is at all possible. For example, suppose we have a function, safely selecting the n -th element of a list:

$$elem : (A : Type) (xs : List\ A) (i : \mathbb{N}) \rightarrow i < length\ xs \rightarrow A$$

If we were to write *elem* without the bounds check $i < length\ xs$, we would get a partial function. Since we can only define total functions in our type theory, we cannot write such a function. However, at run-time, carrying these proofs around makes no sense: type checking has already shown that all calls to *elem* are safe and the proofs do not influence the outcome of *elem*. We want to erase terms of types such as $i < length\ xs$, if we have established that they do not influence the run-time computational behaviour of our functions.

4.1.1 Bove-Capretta method

The *elem* example showed us how we can use propositions to write functions that would otherwise be partial, by asserting properties of the input. The Bove-Capretta method [Bove and Capretta, 2005] generalises this and more: it provides us with a way to transform any (possibly partial) function defined by general recursion into a total, structurally recursive one. The quintessential example of a definition that is not structurally recursive is *quicksort*¹:

$$\begin{aligned} qs &: List\ \mathbb{N} \rightarrow List\ \mathbb{N} \\ qs\ [] &= [] \\ qs\ (x :: xs) &= qs\ (filter\ (gt\ x)\ xs) \uplus x :: qs\ (filter\ (le\ x)\ xs) \end{aligned}$$

The recursive calls are done on $filter\ (gt\ x)\ xs$ and $filter\ (le\ x)\ xs$ instead of just xs , hence *qs* is not structurally recursive. To solve this problem, we create an inductive family describing the call graphs of the original function for every input. Since we can only construct finite values, being able to produce such a call graph essentially means that the function terminates for that input. We can then write a new function that structurally recurses on the call graph. In our quicksort case we get the following inductive family:

$$\begin{aligned} \text{data } qsAcc &: List\ \mathbb{N} \rightarrow Set \text{ where} \\ qsAccNil &: qsAcc\ [] \\ qsAccCons &: (x : \mathbb{N}) (xs : List\ \mathbb{N}) \\ &\quad (h_1 : qsAcc\ (filter\ (gt\ x)\ xs)) \end{aligned}$$

¹In most implementations of functional languages, this definition will not have the same space complexity as the usual in-place version. We are more interested in this function as an example of non-structural recursion and are not too concerned with its complexity.

$$(h_2 : qsAcc (filter (le x) xs)) \\ \rightarrow qsAcc (x :: xs)$$

with the following function definition²

$$\begin{aligned} qs : (xs : List \mathbb{N}) &\rightarrow qsAcc xs \rightarrow List \mathbb{N} \\ qs.nil \quad qsAccNil &= [] \\ qs.cons (qsAccCons x xs h_1 h_2) &= qs (filter (gt x) xs) h_1 ++ \\ &\quad x :: qs (filter (le x) xs) h_2 \end{aligned}$$

Pattern matching on the $qsAcc\ xs$ argument gives us a structurally recursive version of qs . Just as with the *elem* example, we need information from the proof to be able to write this definition in our type theory. In the case of *elem*, we need the proof of $i < length\ xs$ to deal with the (impossible) case where xs is empty. In the qs case, we need $qsAcc\ xs$ to guide the recursion. Even though we actually pattern match on $qsAcc\ xs$ and it therefore seemingly influences the computational behaviour of the function, erasing this argument yields the original qs definition.

4.2 The *Prop* universe in Coq

In Coq we have the *Prop* universe, apart from the *Set* universe. Both universes are base sorts of the hierarchy of sorts, *Type*, i.e. $Prop : Type\ (1)$, $Set : Type\ (1)$ and for every i , $Type\ (i) : Type\ (i + 1)$. As the name suggests, by defining a type to be of sort *Prop*, we “annotate” it to be a logical type, a proposition. Explicitly marking the logical parts like this, makes the development easier to read and understand. More importantly, the extraction mechanism [Letouzey, 2003] now knows what parts are supposed to be logical, hence what parts are to be erased.

In the *sort* example, we would define *isSorted* to be a family of *Props* indexed by $List\ \mathbb{N}$. For the Σ -type, Coq provides two options: *sig* and *ex*, defined as follows:

$$\begin{aligned} \textbf{Inductive } sig\ (A : Type)\ (P : A \rightarrow Prop) : Type &:= \\ \quad exist : \forall x : A, P\ x \rightarrow sig\ P \\ \textbf{Inductive } ex\ (A : Type)\ (P : A \rightarrow Prop) : Prop &:= \\ \quad ex_intro : \forall x : A, P\ x \rightarrow ex\ P \end{aligned}$$

As can be seen above, *sig* differs from *ex* in that the latter is completely logical, whereas *sig* has one informative and one logical field and in its entirety is informative. Since we are interested in the *list* \mathbb{N} part of the Σ -type that is the result type of *sort*, but not the *isSorted* part, we choose the *sig* version.

The extracted version of *sig* consists of a single constructor *exist*, with a single field of type A . Since this is isomorphic to the type A itself, Coq optimises this away

²This definition uses dependent pattern matching [Coquand, 1992], but can be rewritten directly using the elimination operators instead. The important thing here is to notice that we are eliminating the $qsAcc\ xs$ argument.

during extraction. This means $\text{sort} : (xs : \text{List } \mathbb{N}) \rightarrow \Sigma (ys : \text{List } \mathbb{N}) (\text{isSorted } xs \ ys)$ gets extracted to a function $\text{sort}' : \text{List } \mathbb{N} \rightarrow \text{List } \mathbb{N}$.

When erasing all the *Prop* parts from our program, we do want to retain the computational behaviour of the remaining parts. Every function that takes an argument of sort *Prop*, but whose result type is not in *Prop*, needs to be invariant under choice of inhabitant for the *Prop* argument. To force this property, Coq restricts the things we can eliminate a *Prop* into. The general rule is that pattern matching on something of sort *Prop* is allowed if the result type of the function happens to be in *Prop*.

4.2.1 Singleton elimination and homotopy type theory

There are exceptions to this rule: if the argument we are pattern matching on happens to be an *empty* or *singleton definition* of sort *Prop*, we may also eliminate into *Type*. An empty definition is an inductive definition without any constructors. A singleton definition is an inductive definition with precisely one constructor, whose fields are all in *Prop*. Examples of such singleton definitions are conjunction on *Prop* (\wedge) and the accessibility predicate *Acc* used to define functions using well-founded recursion.

Another important example of singleton elimination is elimination on Coq's equality *eq* (where $a = b$ is special notation for $\text{eq } a \ b$), which is defined to be in *Prop*. The inductive family *eq* is defined in the same way as we have defined identity types, hence it is a singleton definition, amenable to singleton elimination. Consider for example the *transport* function:

Definition *transport* : $\forall A, \forall (P : A \rightarrow \text{Type}),$
 $\forall (x \ y : A),$
 $\forall (\text{path} : x = y),$
 $P \ x \rightarrow P \ y .$

Singleton elimination allows us to pattern match on *path* and and eliminate into something of sort *Type*. In the extracted version, the *path* argument gets erased and the $P \ x$ argument is returned. In homotopy type theory, we know that the identity types need not be singletons and can have other inhabitants than just the canonical *refl*, so throwing away the identity proof is not correct. As has been discovered by Michael Shulman³, singleton elimination leads to some sort of inconsistency, if we assume the univalence axiom: we can construct a value $x : \text{bool}$ such that we can prove $x = \text{false}$, even though in the extracted version x normalises to *true*. Assuming univalence, we have two distinct proofs of $\text{bool} = \text{bool}$, namely *refl* and the proof we get from applying univalence to the isomorphism $\text{not} : \text{bool} \rightarrow \text{bool}$. Transporting a value along a path we have obtained from using univalence, is the same as applying the isomorphism. Defining x to be *true* transported along the path obtained from applying univalence to the isomorphism *not*, yields something that is propositionally equal to *false*. If we extract the development, we get a definition of x that ignores the proof of $\text{bool} = \text{bool}$ and just returns *true*.

³<http://homotopytypetheory.org/2012/01/22/univalence-versus-extraction/>

In other words, Coq does not enforce or check proof irrelevance of the types we define to be of sort *Prop*, which internally is fine: it does not allow us to derive falsity using this fact. The extraction mechanism however, does assume that everything admits proof irrelevance. The combination of this along with singleton elimination means that we can prove properties about our programs that no longer hold in the extracted version. It also goes to show that the design decision to define the identity types to be in *Prop* is not compatible with homotopy type theory.

4.2.2 Quicksort example

In the case of *qs* defined using the Bove-Capretta method, we actually want to pattern match on the logical part: *qsAcc xs*. Coq does not allow this if we define the family *qsAcc* to be in *Prop*. However, we can do the pattern matching “manually”, as described in Bertot and Castéran [2004]. We know that we have exactly one inhabitant of *qsAcc xs* for each *xs*, as they represent the call graph of *qs* for the input *xs*, and the pattern matches of the original definition do not overlap, hence each *xs* has a unique call graph. We can therefore easily define and prove the following inversion theorems, that roughly look as follows:

$$\begin{aligned} qsAccInv_0 &: (x : \mathbb{N}) (xs : List\ \mathbb{N}) (qsAcc\ (x :: xs)) \rightarrow qsAcc\ (filter\ (le\ x)\ xs) \\ qsAccInv_1 &: (x : \mathbb{N}) (xs : List\ \mathbb{N}) (qsAcc\ (x :: xs)) \rightarrow qsAcc\ (filter\ (gt\ x)\ xs) \end{aligned}$$

We define the function *qs* just as we originally intended to and add the *qsAcc xs* argument to every pattern match. We then call the inversion theorems for the appropriate recursive calls. Coq still notices that there is a decreasing argument, namely *qsAcc xs*. If we follow this approach, we can define *qsAcc* to be a family in *Prop* and recover the original *qs* definition without the *qsAcc xs* argument using extraction.

In the case of partial functions, we still have to add the missing pattern matches and define impossibility theorems: if we reach that pattern match and we have a proof of our Bove-Capretta predicate for that particular pattern match, we can prove falsity, hence we can use *False_rect* to deal with the missing pattern match.

4.2.3 Impredicativity

So far we have seen how *Prop* differs from *Set* with respect to its restricted elimination rules and its erasure during extraction, but *Prop* has another property that sets it apart from *Set*: *impredicativity*. Impredicativity means that we are able to quantify over something which contains the thing currently being defined. In set theory unrestricted use of this principle leads us to being able to construct Russell’s paradox: the set $R = \{x | x \in x\}$ is an impredicative definition, we quantify over *x*, while we are also defining *x*. Using this definition we can prove that $R \in R$ if and only if $R \notin R$. In type theory, an analogous paradox, Girard’s paradox, arises if we allow for impredicativity via the *Type : Type* rule. However, impredicative definitions are sometimes very useful and benign, in particular when dealing with propositions: we want to be able to write propositions that quantify over propositions, for example:

Definition *demorgan* : $Prop := \forall P Q : Prop,$
 $\sim (P \wedge Q) \rightarrow \sim P \vee \sim Q .$

Coq allows for such definitions as the restrictions on *Prop* prevent us from constructing paradoxes such as Girard's.

4.3 Irrelevance in Agda

In Coq, we put the annotations of something being a proposition in the definition of our inductive type, by defining it to be of sort *Prop*. With Agda's irrelevance mechanism, we instead put the annotations at the places we *use* the proposition, by placing a dot in front of the corresponding type. For example, the type of the *elem* becomes:

elem : (A : Type) (xs : List A) (i : ℕ) → .(i < length xs) → A

We can also mark fields of a record to be irrelevant. In the case of *sort*, we want something similar to the *sig* type from Coq, where second field of the Σ -type is deemed irrelevant. In Agda this can be done as follows:

record Σ_{irr} (A : Type) (B : A → Type) : Type **where**
constructor $\rightarrow, -$
field
fst : A
.snd : B *fst*

To ensure that irrelevant arguments are indeed irrelevant to the computation at hand, Agda has several criteria that it checks. First of all, no pattern matching may be performed on irrelevant arguments, just as is the case with *Prop*. (However, the absurd pattern may be used, if applicable.) Contrary to Coq, singleton elimination is not allowed. Secondly, we need to ascertain that the annotations are preserved: irrelevant arguments may only be passed on to irrelevant contexts. This prevents us from writing a function of type $(A : Type) \rightarrow .A \rightarrow A$.

Another, more important, difference with *Prop* is that irrelevant arguments are ignored by the type checker when checking equality of terms. This can be done safely, even though the terms at hand may in fact be definitionally different, as we never need to appeal to the structure of the value: we cannot pattern match on it. The only thing that we can do with irrelevant arguments is either ignore them or pass them around to other irrelevant contexts.

The reason why the type checker ignoring irrelevant arguments is important, is that it allows us to 'prove' properties about irrelevant arguments in Agda, internally. For example: any function out of an irrelevant type is constant:

irrelevantConstantFunction : {A : Type} {B : Type}
 $\rightarrow (f : .A \rightarrow B) \rightarrow (x\ y : A) \rightarrow f\ x \equiv f\ y$
irrelevantConstantFunction f x y = refl

There is no need to use the congruence rule for \equiv , since the x and y are ignored when the type checker compares $f\ x$ to $f\ y$, when type checking the *refl*. The result can be easily generalised to dependent functions:

$$\begin{aligned} \text{irrelevantConstantDepFunction} &: \{A : \text{Type}\} \{B : A \rightarrow \text{Type}\} \\ &\rightarrow (f : (x : A) \rightarrow B\ x) \rightarrow (x\ y : A) \rightarrow f\ x \equiv f\ y \\ \text{irrelevantConstantDepFunction } f\ x\ y &= \text{refl} \end{aligned}$$

Note that we do not only annotate $(x : A)$ with a dot, but also occurrence of A in the type $B : A \rightarrow \text{Type}$, otherwise we are not allowed to write $B\ x$ as we would use an irrelevant argument in a relevant context. When checking *irrelevantConstantDepFunction*, the term $f\ x \equiv f\ y$ type checks, without having to transport one value along some path, because the types $B\ x$ and $B\ y$ are regarded as definitionally equal by the type checking, ignoring the x and y . Just as before, there is no need to use the (dependent) congruence rule; a *refl* suffices.

We would also like to show that we have proof irrelevance for irrelevant arguments, i.e. we want to prove the following:

$$\text{irrelevantProofIrrelevance} : \{A : \text{Type}\} . (x\ y : A) \rightarrow x \equiv y$$

Agda does not accept this, because the term $x \equiv y$ uses irrelevant arguments in a relevant context: $x \equiv y$. If we instead package the irrelevant arguments in an inductive type, we can prove that the two values of the packaged type are propositionally equal. Consider the following record type with only one irrelevant field:

```
record Squash (A : Type) : Type where
  constructor squash
  field
    .proof : A
```

Using this type, we can now formulate the proof irrelevance principle for irrelevant arguments and prove it:

$$\begin{aligned} \text{squashProofIrrelevance} &: \{A : \text{Type}\} (x\ y : \text{Squash } A) \rightarrow x \equiv y \\ \text{squashProofIrrelevance } x\ y &= \text{refl} \end{aligned}$$

The name “squash type” comes from Nuprl [Constable et al., 1986]: one takes a type and identifies (or “squashes”) all its inhabitants into one unique (up to propositional equality) inhabitant. In homotopy type theory the process of squashing a type is called (-1) -truncation and can also be achieved by defining the following higher inductive type:

```
data (-1)-truncation (A : Type) : Type where
  inhab : A
  all-paths : (x\ y : A) \rightarrow x \equiv y
```

4.3.1 Quicksort example

If we want to mark the *qsAcc xs* argument of the *qs* function as irrelevant, we run into the same problems as we did when we tried to define *qsAcc* as a family in *Prop*: we can no longer pattern match on it. In Coq, we did have a way around this, by using inversion and impossibility theorems to do the pattern matching “manually”. However, if we try such an approach in Agda, its termination checker cannot see that *qsAcc xs* is indeed a decreasing argument and refuses the definition.

4.4 Collapsible families

The approaches we have seen so far let the user indicate what parts of the program are the logical parts and are amenable for erasure. Brady et al. [2004] show that we can let the compiler figure that out by itself instead. The authors propose a series of optimisations for the Epigram system, based on the observation that one often has a lot of redundancy in well-typed terms. If it is the case that one part of a term has to be definitionally equal to another part in order to be well-typed, we can leave out (presuppose) the latter part if we have already established that the term is well-typed.

The authors describe their optimisations in the context of Epigram. In this system, the user writes programs in a high-level language that gets elaborated to programs in a small type theory language. This has the advantage that if we can describe a translation for high-level features, such as dependent pattern matching, to a simple core type theory, the metatheory becomes a lot simpler. The smaller type theory also allows us to specify optimisations more easily, because we do not have to deal with the more intricate, high-level features.

As such, the only things we need to look at, if our goal is to optimise a certain inductive family, are its constructors and its elimination principle. Going back to the *elem* example, we had the *i < length xs* argument. The smaller-than relation can be defined as the following inductive family (in Agda syntax):

```
data _ < _ : ℕ → ℕ → Type where
  ltZ : (y : ℕ)          → Z   < S y
  ltS : (x y : ℕ) → x < y → S x < S y
```

with elimination operator

```
<-elim : (P : (x y : ℕ) → x < y → Type)
        (mZ : (y : ℕ) → P 0 (S y) (ltZ y))
        (mS : (x y : ℕ) → (pf : x < y) → P x y pf → P (S x) (S y) (ltS x y pf))
        (x y : ℕ)
        (pf : x < y)
        → P x y pf
```

and computation rules

$$\begin{aligned}
<-elim\ P\ m_Z\ m_S\ 0\ (S\ y)\ (ltZ\ y) &\mapsto m_Z\ y \\
<-elim\ P\ m_Z\ m_S\ (S\ x)\ (S\ y)\ (ltS\ x\ y\ pf) &\mapsto m_S\ x\ y\ pf\ (<-elim\ P\ m_Z\ m_S\ x\ y\ pf)
\end{aligned}$$

If we look at the computation rules, we see that we can presuppose several things. The first rule has a repeated occurrence of y , so we can presuppose the latter occurrence, the argument of the constructor. In the second rule, the same can be done for x and y . The pf argument can also be erased, as it is never inspected: the only way to inspect pf is via another call the $<-elim$, so by induction it is never inspected. Another thing we observe is that the pattern matches on the indices are disjoint, so we can presuppose the entire target: everything can be recovered from the indices given to the call of $<-elim$.

We have to be careful when making assumptions about values, given their indices. Suppose we have written a function that takes $p : 1 < 1$ as an argument and contains a call to $<-elim$ on p . If we look at the pattern matches on the indices, we may be led to believe that p is of form $ltS\ 0\ 0\ p'$ for some $p' : 0 < 0$ and reduce accordingly. The presupposing only works for *canonical* values, hence we restrict our optimisations to the run-time (evaluation in the empty context), as we know we do not perform reductions under binders in that case and every value is canonical after reduction. The property that every term that is well-typed in the empty context, reduces to a canonical form is called *adequacy* and is a property that is satisfied by Martin-Löf's type theory.

The family $<-elim$ has the property that for indices $x\ y : \mathbb{N}$, its inhabitants $p : x < y$ are uniquely determined by these indices. To be more precise, the following is satisfied: for all $x\ y : \mathbb{N}$, $\vdash p\ q : x < y$ implies $\vdash p \triangleq q$. Families $D : I_0 \rightarrow \dots \rightarrow I_n \rightarrow Type$ such as $<-elim$ are called *collapsible* if they satisfy that for every $i_0 : I_0, \dots, i_n : I_n$, if $\vdash p\ q : D\ i_0\ \dots\ i_n$, then $\vdash p \triangleq q$.

Checking collapsibility of an inductive family is undecidable in general. This can be seen by reducing it to the type inhabitation problem: consider the type $\top + A$. This type is collapsible if and only if A is uninhabited, hence determining with being able to decide collapsibility means we can decide type inhabitation as well. As such, we limit ourselves to a subset that we can recognise, called *concretely collapsible* families. A family $D : I_0 \rightarrow \dots \rightarrow I_n \rightarrow Type$ is concretely collapsible if satisfies the following two properties:

- If we have $\vdash x : D\ i_0\ \dots\ i_n$, for some $i_0 : I_0, \dots, i_n : I_n$, then we can recover its constructor tag by pattern matching on the indices.
- All the non-recursive arguments to the constructors of D can be recovered by pattern matching on the indices.

Note that the first property makes sense because we only have to deal with canonical terms, due to the adequacy property. Checking whether this first property holds can be done by checking whether the indices of the constructors, viewed as patterns, are disjoint. The second property can be checked by pattern matching on the indices of every constructor and checking whether the non-recursive arguments occur as pattern variables.

4.4.1 Erasing concretely collapsible families

If D is a collapsible family, then its elimination operator $D - elim$ is constant in its target, if we fix the indices. This seems to indicate that there might be a possibility to erase the target altogether. Nevertheless, D might have constructors with non-recursive arguments giving us information. Concretely collapsible families satisfy the property that this kind of information can be recovered from the indices, so we can get away with erasing the entire target. Being concretely collapsible means that we have a function at the meta-level (or implementation level) from the indices to the non-recursive, relevant parts of the target. Since this is done by pattern matching on the fully evaluated indices, recovering these parts takes an amount of time that is constant in the size of the given indices. Even though this sounds promising, the complexity of patterns does influence this constant, e.g. the more deeply nested the patterns are, the higher the constant. We now also need the indices to be fully evaluated when eliminating a particular inductive family, whereas that previously might not have been needed. The optimisation is therefore one that gives our dependently typed programs a better space complexity, but not necessarily a better time complexity.

4.4.2 Quicksort example

The accessibility predicates $qsAcc$ form a collapsible family. The pattern matches on the indices in the computation rules for $qsAcc$ are the same pattern matches as those of the original qs definition. There are no overlapping patterns in the original definition, so we can indeed recover the constructor tags from the indices. Also, the non-recursive arguments of $qsAcc$ are precisely those given as indices, hence $qsAcc$ is indeed a (concretely) collapsible family. By the same reasoning, any Bove-Capretta predicate is concretely collapsible, given that the original definition we derived the predicate from, has disjoint pattern matches.

The most important aspect of the collapsibility optimisation is that we have established that everything we need from the value that is to be erased, can be (cheaply) *recovered* from its indices passed to the call to its elimination operator. This means that we have no restrictions on the elimination of collapsible families: we can just write our definition of qs by pattern matching on the $qsAcc$ xs argument. At run-time, the $qsAcc$ xs argument has been erased and the relevant parts are recovered from the indices.

4.5 Internalising collapsibility

Checking whether an inductive family is concretely collapsible is something that can be easily done automatically, as opposed to determining collapsibility in general, which is undecidable. In this section we investigate if we can formulate an internal version of collapsibility, enabling the user to give a proof that a certain family is collapsible, if the compiler fails to notice so itself.

Recall the definition of a collapsible family⁴: given an inductive family D indexed by the type I , D is collapsible if for every index $i : I$ and terms x, y , the following holds:

$$\vdash x, y : D \ i \text{ implies } \vdash x \triangleq y$$

This definition makes use of definitional equality. Since we are working with an intensional type theory, we do not have the *equality reflection rule* at our disposal: there is no rule that tells us that propositional equality implies definitional equality. This might lead us to think that internalising the above definition will not work, as we seemingly cannot say anything about definitional equality from within Martin-Löf's type theory. Let us consider the following variation: for all terms x, y there exists a term p such that

$$\vdash x, y : D \ i \text{ implies } \vdash p : x \equiv y$$

Since Martin-Löf's type theory satisfies the canonicity property, any term p such that $\vdash p : x \equiv y$ reduces to *refl*. The only way for the term to type check, is if $x \triangleq y$, hence in the empty context the equality reflection rule does hold. The converse is also true: definitional equality implies of x and y that $\vdash \text{refl} : x \equiv y$ type checks, hence the latter definition is equal to the original definition of collapsibility.

The variation given above is still not a statement that we can directly prove internally: we need to internalise the implication and replace it by the function space. Doing so yields the following definition: there exists a term p such that:

$$\vdash p : (i : I) \rightarrow (x \ y : D \ i) \rightarrow x \equiv y$$

Or, written as a function in Agda:

$$\begin{aligned} \text{isInternallyCollapsible} &: (I : \text{Type}) (A : I \rightarrow \text{Type}) \rightarrow \text{Type} \\ \text{isInternallyCollapsible } I \ A &= (i : I) \rightarrow (x \ y : A \ i) \rightarrow x \equiv y \end{aligned}$$

We will refer to this definition as *internal collapsibility*. It is easy to see that every internally collapsible family is also collapsible, by canonicity and the fact that *refl* implies definitional equality. However, internally collapsible families do differ from collapsible families as can be seen by considering D to be the family *Id*. By canonicity we have that for any $A : \text{Type}$, $x, y : A$, a term p satisfying $\vdash p : \text{Id } A \ x \ y$ necessarily reduces to *refl*. This means that *Id* is a collapsible family. In contrast, *Id* does not satisfy the internalised condition given above, since this then boils down to the uniqueness of identity proofs principle, which does not hold, as we have discussed.

⁴The definition we originally gave allowed for an arbitrary number of indices. In the following sections we will limit ourselves to the case where we have only one index for presentation purposes. All the results given can be easily generalised to allow more indices.

4.6 Internalising the collapsibility optimisation

In section 4.4.1 we saw how concretely collapsible families can be erased, since all we want to know about the inhabitants can be recovered from its indices. In this section we will try to uncover a similar optimisation for internally collapsible families.

We cannot simply erase the internally collapsible arguments from the function we want to optimise, e.g. given a function $f : (i : I) \rightarrow (x : D\ i) \rightarrow \tau$, we generally cannot produce a function $\tilde{f} : (i : I) \rightarrow \tau$, since we sometimes need the $x : D\ i$ in order for the function to typecheck. However, we can use Agda's irrelevance mechanism to instead generate a function in which the collapsible argument is marked as irrelevant, i.e. we want to write the following function (for the non-dependent case):

```
optimiseFunction :
  (I : Type) (A : I → Type) (B : Type)
  (isInternallyCollapsible I A)
  (f : (i : I) → A i → B)
  → ((i : I) → .(A i) → B)
```

Along with such a function, we should also give a proof that the generated function is equal to the original one in the following sense:

```
optimiseFunctionCorrect :
  (I : Type) (D : I → Type) (B : Type)
  (pf : isInternallyCollapsible I D)
  (f : (i : I) → D i → B)
  (i : I) (x : D i)
  → optimiseFunction I D B pf f i x ≡ f i x
```

If we set out to write the function *optimiseFunction*, after having introduced all the variables, our goal is to produce something of type B . This can be done by using the function f , but then we need a $i : I$ and something of type $D\ i$. We have both, however the $D\ i$ we have is marked as irrelevant, so it may only be passed along to irrelevant contexts, which the function f does not provide, so we cannot use that one. We need to find another way to produce an $D\ i$. We might try to extract it from the proof of *isInternallyCollapsible I D*, but this proof only tells us how the inhabitants of every $D\ i$ are related to each other with propositional equality. From this proof we cannot tell whether some $D\ i$ is inhabited or empty.

The optimisation given for concretely collapsible families need not worry about this. In that case we have a lot more information to work with. We only have to worry about well-typed calls to the elimination operator, so we do not have to deal with deciding whether $D\ i$ is empty or not. Apart from this we only need to recover the non-recursive parts of the erased, canonical term.

If we extend the definition of internal collapsibility with something that decides whether $A\ i$ is empty or not, we get the following definition:

$$\begin{aligned}
& \text{isInternallyCollapsibleDecidable} : (I : \text{Type}) (A : I \rightarrow \text{Type}) \rightarrow \text{Type} \\
& \text{isInternallyCollapsibleDecidable } I \ A = (i : I) \\
& \quad \rightarrow (((x \ y : A \ i) \rightarrow x \equiv y) \otimes (A \ i \oplus A \ i \rightarrow \perp))
\end{aligned}$$

If we then replace the occurrence of *isInternallyCollapsible* in the type signature of *optimiseFunction* with *isInternallyCollapsibleDecidable*

4.6.1 Time complexity issues

Using this definition we do get enough information to write *optimiseFunction*. However, the success of the optimistically named function *optimiseFunction* relies on time complexity the proof given of *isInternallyCollapsibleDecidable* $D \ I$ that is used to recover the erased $A \ i$ value from the index i . In the case of concrete collapsibility this was not that much of an issue, since the way we retrieve the erased values from the indices was constant in the size of the given indices.

Apart from requiring a decision procedure that gives us, for every index $i : I$, an inhabitant of $A \ i$ or a proof that $A \ i$ is empty, we need a bound on the time complexity of this procedure. One approach, taken in Danielsson [2008] to prove time complexities of functions, is to write the functions with a monad that keeps track of how many “ticks” are needed to evaluate the function for the given input, called the *Thunk* monad. $\text{Thunk} : \mathbb{N} \rightarrow \text{Type} \rightarrow \text{Type}$ is implemented as an abstract type that comes with the following primitives:

- $\text{step} : (a : \text{Type}) \rightarrow (n : \mathbb{N}) \rightarrow \text{Thunk } n \ a \rightarrow \text{Thunk } (n + 1) \ a$
- $\text{return} : (a : \text{Type}) \rightarrow (n : \mathbb{N}) \rightarrow a \rightarrow \text{Thunk } n \ a$
- $(\gg) : (a \ b : \text{Type}) \rightarrow (n \ m : \mathbb{N}) \rightarrow \text{Thunk } m \ a \rightarrow (a \rightarrow \text{Thunk } n \ b) \rightarrow \text{Thunk } (m + n) \ b$
- $\text{force} : (a : \text{Type}) \rightarrow (n : \mathbb{N}) \rightarrow \text{Thunk } n \ a$

The user has to write its programs using these primitives. A similar approach has also been used by van Laarhoven⁵ to count the number of comparisons needed for various comparison-based sorting algorithms.

Using this to enforce a time bound on the decision procedure is not too trivial. We first need to establish what kind of time limit we want: do we want a constant time complexity, as we have with the concrete collapsibility optimisation? If we want it to be non-constant, on what variable do we want it to depend?

Apart from these questions, approaches such as the *Thunk* monad, are prone to “cheating”: we can just write our decision procedure the normal way and then write $\text{return } 1 \ \text{decisionProcedure}$ to make sure it has the right type. To prevent this, we can extend the list of primitives in such a way, that the users can write the program completely in this language. Such a language, if it is complete enough, will most likely make writing programs unnecessarily complex for the user.

⁵<http://twanvl.nl/blog/agda/sorting>

Even though we can internalise certain conditions under which certain transformations are safe (preserve definitional equality), along with the transformations, guaranteeing that this transformation actually improves complexity proves to be a lot more difficult.

4.7 Indexed h -propositions and homotopy type theory

In section ?? we have seen that h -propositions are exactly those types that obey proof irrelevance. If we generalise this internal notion to the indexed case we arrive at something we previously have called internal collapsibility. We have also seen that if we restrict ourselves to the empty context, internal collapsibility implies collapsibility. In homotopy type theory, we are interested in postulating extra equalities needed to talk about univalence or higher inductive types. To stress the difference in what contexts we are considering, we will talk about internal collapsible for the empty context case and indexed h -propositions in the other case. In this section we will investigate what these differences mean when trying to optimise our programs.

When postulating extra propositional equalities, we obviously lose the canonicity property, hence we can no longer say that propositional equality implies definitional equality at run-time. The essence of the concrete collapsibility optimisation is that we need not store certain parts of our programs, because we know that they are unique, canonical and can be recovered from other parts of our program. In homotopy type theory we no longer have this canonicity property and may have to make choice in what inhabitant we recover from the indices. As an example of this we will compare two non-indexed types: the unit type and the interval. Both types are h -propositions, so they admit proof irrelevance, but the interval does have two canonical inhabitants that can be distinguished by definitional equality.

```
data I : Set where
  zero : Interval
  one  : Interval
  segment : zero ≡ one
```

The elimination operator for this type is defined in this way:

```
I-elim : (B : I → Type)
  → (b0 : B zero)
  → (b1 : B one)
  → (p : (transport B segment b0) ≡ b1)
  → (i : I) → B i
```

with computation rules⁶:

⁶Apart from giving computation rules for the points, we also need to give a computation rule for the path constructor, *segment*, but as we do not need this rule for the discussion here, we have left it out.

$$\begin{aligned}
I\text{-elim } B \ b_0 \ b_1 \ p \ \text{zero} &\triangleq b_0 \\
I\text{-elim } B \ b_0 \ b_1 \ p \ \text{one} &\triangleq b_1
\end{aligned}$$

In other words, in order to eliminate a value in the interval, we need to tell what has to be done with the endpoints interval and then have to show that this is done in such a way that the path between the endpoints is preserved.

Let us compare the above to the elimination operator for the unit type, \top :

$$\begin{aligned}
\top\text{-elim} : (B : \top \rightarrow \text{Type}) \\
\rightarrow (b : B \ tt) \\
\rightarrow (t : \top) \rightarrow B \ t
\end{aligned}$$

with computation rule:

$$\top\text{-elim } B \ b \ tt \triangleq b$$

If we have canonicity, we can clearly assume every inhabitant of \top to be tt at run-time and erase the t argument from $\top\text{-elim}$. In the case of I , we cannot do this: we have two canonical inhabitants that are propositionally equal, but not definitionally.

Not all is lost, if we consider the non-dependent elimination operator for the interval:

$$\begin{aligned}
I\text{-elim-nondep} : (B : \text{Type}) \\
\rightarrow (b_0 : B) \\
\rightarrow (b_1 : B) \\
\rightarrow (p : b_0 \equiv b_1) \\
\rightarrow I \rightarrow B
\end{aligned}$$

then it is easy to see that all such functions are constant functions, with respect to propositional equality. If we erase the I argument and presuppose it to be zero , we will get a new function that is propositionally equal to the original one. However, it is definitional equality that we are after. We can define the following two functions:

$$\begin{aligned}
I\text{-id} : I \rightarrow I \\
I\text{-id} &= I\text{-elim-nondep } I \ \text{zero} \ \text{one} \ \text{segment} \\
I\text{-const-zero} : I \rightarrow I \\
I\text{-const-zero} &= I\text{-elim-nondep } I \ \text{zero} \ \text{zero} \ \text{refl}
\end{aligned}$$

If we presuppose and erase the I argument to be zero in the $I\text{-id}$ case, we would get definitionally different behaviour. In the case of $I\text{-const-zero}$, it does not matter if we presuppose the argument to be zero or one , since this function is also definitionally constant. This is because for the refl to type check, b_0 and b_1 have to be definitionally equal. So if we want to optimise the elimination operators of higher inductive types that are h -propositions, such as the interval, we need to look at what paths the non-trivial paths are mapped to. If these are all mapped

to *refl*, then the points all get mapped to definitionally equal points. Checking such a property can become difficult, as we can tell from this rather silly example:

```
data ℕ-truncated : Type where
  0 : ℕ-truncated
  S : (n : ℕ-truncated) → ℕ-truncated
  equalTo0 : (n : ℕ-truncated) → 0 ≡ n
```

with non-dependent eliminator:

```
ℕ-truncated-elim-nondep : (B : Type)
  → (b₀ : B)
  → (b_S : B → B)
  → (p : (b : B) → b₀ ≡ b)
  → ℕ-truncated → B
```

If we were to check that all paths between 0 and *n* are mapped to a *refl* between inhabitants of *B*, we have to check that *p* satisfies this property, which we cannot do.

4.7.1 Internally optimising *h*-propositions

The optimisation given in section 4.6 of course still is a valid transformation for the homotopy type theory case. The proof of a family $D : I \rightarrow \text{Type}$ being an indexed *h*-proposition is again not enough for us to be able to write the *optimiseFunction* term. What we called *isInternallyCollapsibleDecidable* is that we internally need a witness of the fact that every *h*-proposition in the family is either contractible or empty, so we could have written the property as follows:

```
isIndexedhPropDecidable : (I : Type) (A : I → Type) → Type
isIndexedhPropDecidable I A = (i : I)
  → (isContractible (A i)) ⊕ (A i → ⊥)
```

4.8 Conclusions

In this chapter we have looked at various ways of dealing with types that are purely logical, called propositions. Coq and Agda both provide mechanisms to in a way “truncate” a type into a proposition. The first takes this approach by allowing the user to annotate a type as being a proposition when defining the type. Making sure it is a proposition and has no computational effect on non-propositions is handled by limiting the elimination of these propositions: we may only eliminate into other propositions. Singleton elimination is an exception to this rule, which does not play well with homotopy type theory and the univalence axiom. Proof irrelevance of the propositions in Coq is assumed when extracting a development, but not something that is enforced inside Coq, nor is it provable

internally. Using univalence we can construct a term that behaves differently in Coq as it does in the extracted version.

Agda allows the user to indicate that a type is a proposition when referring to that type, instead of having to annotate it when defining it. Agda enforces the proof irrelevance by ensuring that inhabitants of an annotated type are never scrutinised in a pattern match and may only be passed onto other irrelevant contexts. In contrast to Coq's mechanism, it does not allow for singleton elimination, but unlike Coq, it does enable the user to prove properties of the annotated types in Agda itself. As such, we can construct a squash type that is isomorphic to the (-1) -truncation from homotopy type theory, defined as a higher inductive type.

Instead of truncating a type such that it becomes a proposition, we can also let the compiler recognise whether a type is a proposition or not. This is the approach that the collapsible families optimisation takes in Epigram. The definition of collapsibility is reminiscent of the definition of h -proposition, albeit it is an indexed version that uses definitional equality instead of propositional equality. The optimisation specifically focuses on families of propositions.

Recognising whether an inductive family is a collapsible family is undecidable, so the actual optimisation restricts itself to a subset called concretely collapsible families. To improve on this, we internalise the notion of collapsibility, allowing the user to provide a proof if the compiler fails to notice this property. We show that this notion of internal collapsibility is a subset of collapsibility. We also try to internalise the optimisation, but since the time complexity of the optimised function heavily depends on the user-provided proof, we cannot be sure whether it the "optimised" version actually improves on the complexity. We have looked at ways to enforce time complexities in the user-provided proofs. Our conclusion is that this is not viable.

As we have mentioned previously, collapsible families look a lot like families of h -propositions. When internalising the collapsibility concept and the optimisation, we only considered the non-homotopy type theory case, i.e. no univalence and no higher inductive types. We have looked at extending the optimisations to the homotopy type theory case, but as we lose canonicity the optimised versions may no longer yield the same results as the original function, with respect to definitional equality. We have identified cases in which this is the case and cases in which definitional equality actually is preserved. We also argue that detecting whether the latter is the case, is not tractable.

Chapter 5

Discussion

One of the main goals of this project was to establish whether homotopy type theory is an interesting language to do dependently typed programming in. As it is incompatible with dependent pattern matching in general, it seems like we are taking a step backwards. However, univalence and higher inductive types can become the two steps forward. Univalence means that we can transport definitions along isomorphisms, which saves us a great deal of writing boring code applying the *to* and *from* parts of the isomorphisms in the right places. It also implies function extensionality, which is indispensable when proving properties about programs.

We have also seen the usefulness of higher inductive types. They allow us to define quotient types. It is all too easy to come up with a higher inductive type that has more structure than is desired: one quickly runs into *coherence issues*: the resulting type has too much different equalities at higher levels than is needed. The original definition of quotient types also suffered from this issue: we want it to be a *h*-set, but as could be seen from a simple example, one could easily define the circle: the simplest type that is not a *h*-set. Therefore one usually needs to truncate the higher inductive type to a certain level, e.g. take the 0-truncation in the case of quotients. Truncating a type does mean that we have extra conditions that we need to satisfy when eliminating something of that particular type. In a programming setting, one typically only encounters *h*-sets, except for univalent universes of *h*-sets. Eliminating into a *h*-set means that the extra conditions stemming from 0-truncation are automatically satisfied, so in programming this need not be too much of a problem.

For these two steps to be actual steps forward, there is still a lot of work that needs to be done. The most obvious and possibly most difficult problem is determining the computational content of the univalence axiom. For programming purposes however, having a type theory in which everything is 1-truncated is already a big improvement, since, as we have already mentioned, most types in programming applications are *h*-sets.

Giving up pattern matching altogether is quite drastic. There are still a lot of cases in which (dependent) pattern matching is still valid and can be transformed to an expression using only elimination principles. An interesting future research di-

rection is to take the elaboration process described in [Goguen et al. \[2006\]](#), which critically depends on axiom K , and see how one can uncover conditions in which K is not necessary for the elaboration to work.

There is also a lot of work to be done on higher inductive types. As of yet, a well-defined syntax for higher inductive types and a generic way to derive the induction principles is lacking. It has also been noted [[Lumsdaine, 2012](#)] that every higher inductive type that has higher path constructors in its definition, can be rewritten to an equivalent form that only has path constructors that construct paths between points (a so called 1-HIT). Having a mechanism that automatically translates the definition of a higher inductive type to a 1-HIT, also means that we only have to care about these cases when devising induction principles. Having a form of pattern matching for higher inductive types is also a research direction that can help make higher inductive types significantly more easy to work with.

In ??, we have seen that in traditional Martin-Löf's type theory, propositional equality coincides with definitional equality at "run-time" (i.e. in the empty context). This property makes it possible to internalise optimisations: one could create a system in which we provide rules akin to the GHC REWRITE rules, but along with a proof of correctness. In homotopy type theory, we also want to have non-canonical proofs of propositional equality at run-time, so we lose this property. A further investigation of when propositional equality still does imply definitional equality might be an interesting research direction. Another interesting thing to look at is the question whether we really need definitional equality, i.e. identify cases in which we can safely replace something by something else that is propositionally but not necessarily definitionally equal.

Coming back to the main research question:

What is homotopy type theory and why is it interesting to do programming in?

There is evidence that homotopy type theory is an interesting language to program in, but there are still a lot of things that need to be worked out and solved before it becomes useful.

Acknowledgements

Front and foremost, I would like to thank my supervisor Wouter for the guidance writing this thesis. Secondly, Paul deserves a mention for hosting me in Bordeaux, renewing my motivation to work on this project and letting me experience French culture, including but not limited to air traffic control and train strikes. The people from BBL682 and BBL681 provided much needed distraction. Last but not least, I would like to thank my parents for their unconditional support.

Bibliography

- T. Altenkirch, T. Anberrée, and N. Li. Definable quotients in type theory.
- Y. Bertot and P. Castéran. *Interactive theorem proving and program development: Coq'Art: the calculus of inductive constructions*. Springer-Verlag New York Incorporated, 2004.
- A. Bove and V. Capretta. Modelling general recursion in type theory. *Mathematical Structures in Computer Science*, 15(4):671–708, 2005.
- E. Brady, C. McBride, and J. McKinna. Inductive families need not store their indices. In *Types for Proofs and Programs*, pages 115–129. Springer, 2004.
- R. Constable, S. Allen, H. Bromley, W. Cleaveland, J. Cremer, R. Harper, D. Howe, T. Knoblock, N. Mendler, P. Panangaden, et al. *Implementing mathematics with the Nuprl proof development system*. Prentice Hall, 1986.
- T. Coquand. Pattern matching with dependent types. In *Informal proceedings of Logical Frameworks*, volume 92, pages 66–79, 1992.
- N. A. Danielsson. Lightweight semiformal time complexity analysis for purely functional data structures. In *ACM SIGPLAN Notices*, volume 43, pages 133–144. ACM, 2008.
- N. A. Danielsson. Postulated computing quotients are unsound. online, <https://lists.chalmers.se/pipermail/agda/2012/004052.html>, 2012. [Agda mailing list post].
- H. Goguen, C. McBride, and J. McKinna. Eliminating dependent pattern matching. In *Algebra, Meaning, and Computation*, pages 521–540. Springer, 2006.
- M. Hofmann and T. Streicher. The groupoid interpretation of type theory. In *In Venice Festschrift*, 1996.
- N. Kraus and C. Sattler. Universe n is not an n -type. online, <http://homotopytypetheory.org/2013/05/15/universe-n-is-not-an-n-type/>, 2013. [blog post].
- P. Letouzey. A new extraction for Coq. In *Types for proofs and programs*, pages 200–219. Springer, 2003.
- D. R. Licata. Running Circles Around (In) Your Proof Assistant; or, Quotients that Compute. online, <http://homotopytypetheory.org/2011/04/23/running-circles-around-in-your-proof-assistant/>, 2011. [blog post].

- D. R. Licata. Running Circles Around (In) Your Proof Assistant; or, Quotients that Compute. online, <http://homotopytypetheory.org/2012/11/12/abstract-types-with-isomorphic-types/>, 2012. [blog post].
- D. R. Licata and R. Harper. Canonicity for 2-dimensional type theory. In *ACM SIGPLAN Notices*, volume 47, pages 337–348. ACM, 2012.
- D. R. Licata and M. Shulman. Calculating the fundamental group of the circle in homotopy type theory. *arXiv preprint arXiv:1301.3443*, 2013.
- P. L. Lumsdaine. Reducing all HIT’s to 1-HIT’s. online, <http://homotopytypetheory.org/2012/05/07/reducing-all-hits-to-1-hits/>, 2012. [blog post].
- P. Martin-Löf. Constructive mathematics and computer programming. In *Proc. of a discussion meeting of the Royal Society of London on Mathematical logic and programming languages*, pages 167–184. Prentice-Hall, Inc., 1985.
- J. C. Mitchell and G. D. Plotkin. Abstract types have existential type. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 10(3):470–502, 1988.
- M. Sozeau, N. Tabareau, et al. Univalence for free. 2013.
- V. Voevodsky. Univalent foundations. online, http://www.math.ias.edu/~vladimir/Site3/Univalent_Foundations_files/2011_UPenn.pdf, 2011. [presentation at University of Pennsylvania].
- P. Wadler. Views: A way for pattern matching to cohabit with data abstraction. In *Proceedings of the 14th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 307–313. ACM, 1987.

Index of symbols

asdf

propositional equality

Index

continuous deformation, 4

definable quotient, 21
dictionary, 20

equivalence, 16
extensional type theory, 5

intensional type theory, 5

path space, 4
propositional equality, 5

quotient, 19

uniqueness of identity proofs, 7