

Homotopy Type Theory

Univalent Foundations of Mathematics

The Univalent Foundations Program

Institute for Advanced Study

“Homotopy Type Theory: Univalent Foundations of Mathematics”

© 2013 The Univalent Foundations Program

Book version: first-edition-1174-g29279f5

MSC 2010 classification: 03-02, 55-02, 03B15

This work is licensed under the *Creative Commons Attribution-ShareAlike 3.0 Unported License*. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/>.

This book is freely available at <http://homotopytypetheory.org/book/>.

Acknowledgment

Apart from the generous support from the Institute for Advanced Study, some contributors to the book were partially or fully supported by the following agencies and grants:

- Association of Members of the Institute for Advanced Study: a grant to the Institute for Advanced Study
- Agencija za raziskovalno dejavnost Republike Slovenije: P1-0294, N1-0011.
- Air Force Office of Scientific Research: FA9550-11-1-0143, and FA9550-12-1-0370.

This material is based in part upon work supported by the AFOSR under the above awards. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the AFOSR.

- Engineering and Physical Sciences Research Council: EP/G034109/1, EP/G03298X/1.
- European Union’s 7th Framework Programme under grant agreement nr. 243847 (ForMath).
- National Science Foundation: DMS-1001191, DMS-1100938, CCF-1116703, and DMS-1128155.

This material is based in part upon work supported by the National Science Foundation under the above awards. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

- The Simonyi Fund: a grant to the Institute for Advanced Study

Preface

IAS Special Year on Univalent Foundations

A Special Year on Univalent Foundations of Mathematics was held in 2012-13 at the Institute for Advanced Study, School of Mathematics, organized by Steve Awodey, Thierry Coquand, and Vladimir Voevodsky. The following people were the official participants.

Peter Aczel	Eric Finster	Alvaro Pelayo
Benedikt Ahrens	Daniel Grayson	Andrew Polonsky
Thorsten Altenkirch	Hugo Herbelin	Michael Shulman
Steve Awodey	André Joyal	Matthieu Sozeau
Bruno Barras	Dan Licata	Bas Spitters
Andrej Bauer	Peter Lumsdaine	Benno van den Berg
Yves Bertot	Assia Mahboubi	Vladimir Voevodsky
Marc Bezem	Per Martin-Löf	Michael Warren
Thierry Coquand	Sergey Melikhov	Noam Zeilberger

There were also the following students, whose participation was no less valuable.

Carlo Angiuli	Guillaume Brunerie	Egbert Rijke
Anthony Bordg	Chris Kapulkin	Kristina Sojakova

In addition, there were the following short- and long-term visitors, including student visitors, whose contributions to the Special Year were also essential.

Jeremy Avigad	Richard Garner	Nuo Li
Cyril Cohen	Georges Gonthier	Zhaohui Luo
Robert Constable	Thomas Hales	Michael Nahas
Pierre-Louis Curien	Robert Harper	Erik Palmgren
Peter Dybjer	Martin Hofmann	Emily Riehl
Martín Escardó	Pieter Hofstra	Dana Scott
Kuen-Bang Hou	Joachim Kock	Philip Scott
Nicola Gambino	Nicolai Kraus	Sergei Soloviev

About this book

We did not set out to write a book. The present work has its origins in our collective attempts to develop a new style of “informal type theory” that can be read and understood by a human being, as a complement to a formal proof that can be checked by a machine. Univalent foundations is closely tied to the idea of a foundation of mathematics that can be implemented in a computer proof assistant. Although such a formalization is not part of this book, much of the material presented here was actually done first in the fully formalized setting inside a proof assistant, and only later “unformalized” to arrive at the presentation you find before you — a remarkable inversion of the usual state of affairs in formalized mathematics.

Each of the above-named individuals contributed something to the Special Year — and so to this book — in the form of ideas, words, or deeds. The spirit of collaboration that prevailed throughout the year was truly extraordinary.

Special thanks are due to the Institute for Advanced Study, without which this book would obviously never have come to be. It proved to be an ideal setting for the creation of this new branch of mathematics: stimulating, congenial, and supportive. May some trace of this unique atmosphere linger in the pages of this book, and in the future development of this new field of study.

The Univalent Foundations Program
Institute for Advanced Study
Princeton, April 2013

Contents

Introduction	1
I Foundations	15
1 Type theory	17
1.1 Type theory versus set theory	17
1.2 Function types	21
1.3 Universes and families	23
1.4 Dependent function types (Π -types)	24
1.5 Product types	26
1.6 Dependent pair types (Σ -types)	29
1.7 Coproduct types	32
1.8 The type of booleans	33
1.9 The natural numbers	35
1.10 Pattern matching and recursion	38
1.11 Propositions as types	39
1.12 Identity types	45
Notes	51
Exercises	53
2 Homotopy type theory	55
2.1 Types are higher groupoids	58
2.2 Functions are functors	66
2.3 Type families are fibrations	67
2.4 Homotopies and equivalences	70
2.5 The higher groupoid structure of type formers	74
2.6 Cartesian product types	75
2.7 Σ -types	77
2.8 The unit type	80
2.9 Π -types and the function extensionality axiom	80
2.10 Universes and the univalence axiom	83
2.11 Identity type	84
2.12 Coproducts	86
2.13 Natural numbers	88
2.14 Example: equality of structures	90

2.15 Universal properties	93
Notes	95
Exercises	97
3 Sets and logic	99
3.1 Sets and n -types	99
3.2 Propositions as types?	101
3.3 Mere propositions	103
3.4 Classical vs. intuitionistic logic	104
3.5 Subsets and propositional resizing	106
3.6 The logic of mere propositions	107
3.7 Propositional truncation	108
3.8 The axiom of choice	110
3.9 The principle of unique choice	111
3.10 When are propositions truncated?	112
3.11 Contractibility	114
Notes	117
Exercises	117
4 Equivalences	121
4.1 Quasi-inverses	122
4.2 Half adjoint equivalences	124
4.3 Bi-invertible maps	128
4.4 Contractible fibers	128
4.5 On the definition of equivalences	129
4.6 Surjections and embeddings	130
4.7 Closure properties of equivalences	131
4.8 The object classifier	133
4.9 Univalence implies function extensionality	135
Notes	137
Exercises	137
5 Induction	139
5.1 Introduction to inductive types	139
5.2 Uniqueness of inductive types	141
5.3 W-types	144
5.4 Inductive types are initial algebras	147
5.5 Homotopy-inductive types	149
5.6 The general syntax of inductive definitions	153
5.7 Generalizations of inductive types	156
5.8 Identity types and identity systems	159
Notes	162
Exercises	163

6 Higher inductive types	167
6.1 Introduction	167
6.2 Induction principles and dependent paths	169
6.3 The interval	173
6.4 Circles and spheres	174
6.5 Suspensions	176
6.6 Cell complexes	179
6.7 Hubs and spokes	180
6.8 Pushouts	181
6.9 Truncations	185
6.10 Quotients	187
6.11 Algebra	192
6.12 The flattening lemma	196
6.13 The general syntax of higher inductive definitions	201
Notes	203
Exercises	204
7 Homotopy n-types	205
7.1 Definition of n -types	205
7.2 Uniqueness of identity proofs and Hedberg's theorem	208
7.3 Truncations	212
7.4 Colimits of n -types	217
7.5 Connectedness	221
7.6 Orthogonal factorization	225
7.7 Modalities	230
Notes	233
Exercises	234
Bibliography	239
Index of symbols	245
Index	251

Introduction

Homotopy type theory is a new branch of mathematics that combines aspects of several different fields in a surprising way. It is based on a recently discovered connection between *homotopy theory* and *type theory*. Homotopy theory is an outgrowth of algebraic topology and homological algebra, with relationships to higher category theory; while type theory is a branch of mathematical logic and theoretical computer science. Although the connections between the two are currently the focus of intense investigation, it is increasingly clear that they are just the beginning of a subject that will take more time and more hard work to fully understand. It touches on topics as seemingly distant as the homotopy groups of spheres, the algorithms for type checking, and the definition of weak ∞ -groupoids.

Homotopy type theory also brings new ideas into the very foundation of mathematics. On the one hand, there is Voevodsky’s subtle and beautiful *univalence axiom*. The univalence axiom implies, in particular, that isomorphic structures can be identified, a principle that mathematicians have been happily using on workdays, despite its incompatibility with the “official” doctrines of conventional foundations. On the other hand, we have *higher inductive types*, which provide direct, logical descriptions of some of the basic spaces and constructions of homotopy theory: spheres, cylinders, truncations, localizations, etc. Both ideas are impossible to capture directly in classical set-theoretic foundations, but when combined in homotopy type theory, they permit an entirely new kind of “logic of homotopy types”.

This suggests a new conception of foundations of mathematics, with intrinsic homotopical content, an “invariant” conception of the objects of mathematics — and convenient machine implementations, which can serve as a practical aid to the working mathematician. This is the *Univalent Foundations* program. The present book is intended as a first systematic exposition of the basics of univalent foundations, and a collection of examples of this new style of reasoning — but without requiring the reader to know or learn any formal logic, or to use any computer proof assistant.

We emphasize that homotopy type theory is a young field, and univalent foundations is very much a work in progress. This book should be regarded as a “snapshot” of just one portion of the field, taken at the time it was written, rather than a polished exposition of a completed edifice. As we will discuss briefly later, there are many aspects of homotopy type theory that are not yet fully understood — and some that are not even touched upon here. The ultimate theory will almost certainly not look exactly like the one described in this book, but it will surely be *at least* as capable and powerful; we therefore believe that univalent foundations will eventually become a viable alternative to set theory as the “implicit foundation” for the unformalized mathematics done by most mathematicians.

Type theory

Type theory was originally invented by Bertrand Russell [Rus08], as a device for blocking the paradoxes in the logical foundations of mathematics that were under investigation at the time. It was developed further by many people over the next few decades, particularly Church [Chu40, Chu41] who combined it with his λ -calculus. Although it is not generally regarded as the foundation for classical mathematics, set theory being more customary, type theory still has numerous applications, especially in computer science and the theory of programming languages [Pie02]. Per Martin-Löf [ML98, ML75, ML82, ML84], among others, developed a “predicative” modification of Church’s type system, which is now usually called dependent, constructive, intuitionistic, or simply *Martin-Löf type theory*. This is the basis of the system that we consider here; it was originally intended as a rigorous framework for the formalization of constructive mathematics. In what follows, we will often use “type theory” to refer specifically to this system and similar ones, although type theory as a subject is much broader (see [Som10, KLN04] for the history of type theory).

In type theory, unlike set theory, objects are classified using a primitive notion of *type*, similar to the data-types used in programming languages. These elaborately structured types can be used to express detailed specifications of the objects classified, giving rise to principles of reasoning about these objects. To take a very simple example, the objects of a product type $A \times B$ are known to be of the form (a, b) , and so one automatically knows how to construct them and how to decompose them. Similarly, an object of function type $A \rightarrow B$ can be acquired from an object of type B parametrized by objects of type A , and can be evaluated at an argument of type A . This rigidly predictable behavior of all objects (as opposed to set theory’s more liberal formation principles, allowing inhomogeneous sets) is one aspect of type theory that has led to its extensive use in verifying the correctness of computer programs. The clear reasoning principles associated with the construction of types also form the basis of modern *computer proof assistants*, which are used for formalizing mathematics and verifying the correctness of formalized proofs. We return to this aspect of type theory below.

One problem in understanding type theory from a mathematical point of view, however, has always been that the basic concept of *type* is unlike that of *set* in ways that have been hard to make precise. We believe that the new idea of regarding types, not as strange sets (perhaps constructed without using classical logic), but as spaces, viewed from the perspective of homotopy theory, is a significant step forward. In particular, it solves the problem of understanding how the notion of equality of elements of a type differs from that of elements of a set.

In homotopy theory one is concerned with spaces and continuous mappings between them, up to homotopy. A *homotopy* between a pair of continuous maps $f : X \rightarrow Y$ and $g : X \rightarrow Y$ is a continuous map $H : X \times [0, 1] \rightarrow Y$ satisfying $H(x, 0) = f(x)$ and $H(x, 1) = g(x)$. The homotopy H may be thought of as a “continuous deformation” of f into g . The spaces X and Y are said to be *homotopy equivalent*, $X \simeq Y$, if there are continuous maps going back and forth, the composites of which are homotopical to the respective identity mappings, i.e., if they are isomorphic “up to homotopy”. Homotopy equivalent spaces have the same algebraic invariants (e.g., homology, or the fundamental group), and are said to have the same *homotopy type*.

Homotopy type theory

Homotopy type theory (HoTT) interprets type theory from a homotopical perspective. In homotopy type theory, we regard the types as “spaces” (as studied in homotopy theory) or higher groupoids, and the logical constructions (such as the product $A \times B$) as homotopy-invariant con-

structions on these spaces. In this way, we are able to manipulate spaces directly without first having to develop point-set topology (or any combinatorial replacement for it, such as the theory of simplicial sets). To briefly explain this perspective, consider first the basic concept of type theory, namely that the *term* a is of *type* A , which is written:

$$a : A.$$

This expression is traditionally thought of as akin to:

$$\text{"}a\text{ is an element of the set }A\text{"}.$$

However, in homotopy type theory we think of it instead as:

$$\text{"}a\text{ is a point of the space }A\text{"}.$$

Similarly, every function $f : A \rightarrow B$ in type theory is regarded as a continuous map from the space A to the space B .

We should stress that these “spaces” are treated purely homotopically, not topologically. For instance, there is no notion of “open subset” of a type or of “convergence” of a sequence of elements of a type. We only have “homotopical” notions, such as paths between points and homotopies between paths, which also make sense in other models of homotopy theory (such as simplicial sets). Thus, it would be more accurate to say that we treat types as ∞ -groupoids; this is a name for the “invariant objects” of homotopy theory which can be presented by topological spaces, simplicial sets, or any other model for homotopy theory. However, it is convenient to sometimes use topological words such as “space” and “path”, as long as we remember that other topological concepts are not applicable.

(It is tempting to also use the phrase *homotopy type* for these objects, suggesting the dual interpretation of “a type (as in type theory) viewed homotopically” and “a space considered from the point of view of homotopy theory”. The latter is a bit different from the classical meaning of “homotopy type” as an *equivalence class* of spaces modulo homotopy equivalence, although it does preserve the meaning of phrases such as “these two spaces have the same homotopy type”.)

The idea of interpreting types as structured objects, rather than sets, has a long pedigree, and is known to clarify various mysterious aspects of type theory. For instance, interpreting types as sheaves helps explain the intuitionistic nature of type-theoretic logic, while interpreting them as partial equivalence relations or “domains” helps explain its computational aspects. It also implies that we can use type-theoretic reasoning to study the structured objects, leading to the rich field of categorical logic. The homotopical interpretation fits this same pattern: it clarifies the nature of *identity* (or equality) in type theory, and allows us to use type-theoretic reasoning in the study of homotopy theory.

The key new idea of the homotopy interpretation is that the logical notion of identity $a = b$ of two objects $a, b : A$ of the same type A can be understood as the existence of a path $p : a \sim b$ from point a to point b in the space A . This also means that two functions $f, g : A \rightarrow B$ can be identified if they are homotopic, since a homotopy is just a (continuous) family of paths $p_x : f(x) \sim g(x)$ in B , one for each $x : A$. In type theory, for every type A there is a (formerly somewhat mysterious) type Id_A of identifications of two objects of A ; in homotopy type theory, this is just the *path space* A^I of all continuous maps $I \rightarrow A$ from the unit interval. In this way, a term $p : \text{Id}_A(a, b)$ represents a path $p : a \sim b$ in A .

The idea of homotopy type theory arose around 2006 in independent work by Awodey and Warren [AW09] and Voevodsky [Voe06], but it was inspired by Hofmann and Streicher’s earlier

groupoid interpretation [HS98]. Indeed, higher-dimensional category theory (particularly the theory of weak ∞ -groupoids) is now known to be intimately connected to homotopy theory, as proposed by Grothendieck and now being studied intensely by mathematicians of both sorts. The original semantic models of Awodey–Warren and Voevodsky use well-known notions and techniques from homotopy theory which are now also in use in higher category theory, such as Quillen model categories and Kan simplicial sets.

In particular, Voevodsky constructed an interpretation of type theory in Kan simplicial sets, and recognized that this interpretation satisfied a further crucial property which he dubbed *univalence*. This had not previously been considered in type theory (although Church’s principle of extensionality for propositions turns out to be a very special case of it, and Hofmann and Streicher had considered another special case under the name “universe extensionality”). Adding univalence to type theory in the form of a new axiom has far-reaching consequences, many of which are natural, simplifying and compelling. The univalence axiom also further strengthens the homotopical view of type theory, since it holds in the simplicial model and other related models, while failing under the view of types as sets.

Univalent foundations

Very briefly, the basic idea of the univalence axiom can be explained as follows. In type theory, one can have a universe \mathcal{U} , the terms of which are themselves types, $A : \mathcal{U}$, etc. Those types that are terms of \mathcal{U} are commonly called *small* types. Like any type, \mathcal{U} has an identity type $\text{Id}_{\mathcal{U}}$, which expresses the identity relation $A = B$ between small types. Thinking of types as spaces, \mathcal{U} is a space, the points of which are spaces; to understand its identity type, we must ask, what is a path $p : A \sim B$ between spaces in \mathcal{U} ? The univalence axiom says that such paths correspond to homotopy equivalences $A \simeq B$, (roughly) as explained above. A bit more precisely, given any (small) types A and B , in addition to the primitive type $\text{Id}_{\mathcal{U}}(A, B)$ of identifications of A with B , there is the defined type $\text{Equiv}(A, B)$ of equivalences from A to B . Since the identity map on any object is an equivalence, there is a canonical map,

$$\text{Id}_{\mathcal{U}}(A, B) \rightarrow \text{Equiv}(A, B).$$

The univalence axiom states that this map is itself an equivalence. At the risk of oversimplifying, we can state this succinctly as follows:

Univalence Axiom: $(A = B) \simeq (A \simeq B)$.

In other words, identity is equivalent to equivalence. In particular, one may say that “equivalent types are identical”. However, this phrase is somewhat misleading, since it may sound like a sort of “skeletality” condition which *collapses* the notion of equivalence to coincide with identity, whereas in fact univalence is about *expanding* the notion of identity so as to coincide with the (unchanged) notion of equivalence.

From the homotopical point of view, univalence implies that spaces of the same homotopy type are connected by a path in the universe \mathcal{U} , in accord with the intuition of a classifying space for (small) spaces. From the logical point of view, however, it is a radically new idea: it says that isomorphic things can be identified! Mathematicians are of course used to identifying isomorphic structures in practice, but they generally do so by “abuse of notation”, or some other informal device, knowing that the objects involved are not “really” identical. But in this new foundational scheme, such structures can be formally identified, in the logical sense that every

property or construction involving one also applies to the other. Indeed, the identification is now made explicit, and properties and constructions can be systematically transported along it. Moreover, the different ways in which such identifications may be made themselves form a structure that one can (and should!) take into account.

Thus in sum, for points A and B of the universe \mathcal{U} (i.e., small types), the univalence axiom identifies the following three notions:

- (logical) an identification $p : A = B$ of A and B
- (topological) a path $p : A \rightsquigarrow B$ from A to B in \mathcal{U}
- (homotopical) an equivalence $p : A \simeq B$ between A and B .

Higher inductive types

One of the classical advantages of type theory is its simple and effective techniques for working with inductively defined structures. The simplest nontrivial inductively defined structure is the natural numbers, which is inductively generated by zero and the successor function. From this statement one can algorithmically extract the principle of mathematical induction, which characterizes the natural numbers. More general inductive definitions encompass lists and well-founded trees of all sorts, each of which is characterized by a corresponding “induction principle”. This includes most data structures used in certain programming languages; hence the usefulness of type theory in formal reasoning about the latter. If conceived in a very general sense, inductive definitions also include examples such as a disjoint union $A + B$, which may be regarded as “inductively” generated by the two injections $A \rightarrow A + B$ and $B \rightarrow A + B$. The “induction principle” in this case is “proof by case analysis”, which characterizes the disjoint union.

In homotopy theory, it is natural to consider also “inductively defined spaces” which are generated not merely by a collection of *points*, but also by collections of *paths* and higher paths. Classically, such spaces are called *CW complexes*. For instance, the circle S^1 is generated by a single point and a single path from that point to itself. Similarly, the 2-sphere S^2 is generated by a single point b and a single two-dimensional path from the constant path at b to itself, while the torus T^2 is generated by a single point, two paths p and q from that point to itself, and a two-dimensional path from $p \cdot q$ to $q \cdot p$.

By using the identification of paths with identities in homotopy type theory, these sort of “inductively defined spaces” can be characterized in type theory by “induction principles”, entirely analogously to classical examples such as the natural numbers and the disjoint union. The resulting *higher inductive types* give a direct “logical” way to reason about familiar spaces such as spheres, which (in combination with univalence) can be used to perform familiar arguments from homotopy theory, such as calculating homotopy groups of spheres, in a purely formal way. The resulting proofs are a marriage of classical homotopy-theoretic ideas with classical type-theoretic ones, yielding new insight into both disciplines.

Moreover, this is only the tip of the iceberg: many abstract constructions from homotopy theory, such as homotopy colimits, suspensions, Postnikov towers, localization, completion, and spectrification, can also be expressed as higher inductive types. Many of these are classically constructed using Quillen’s “small object argument”, which can be regarded as a finite way of algorithmically describing an infinite CW complex presentation of a space, just as “zero and successor” is a finite algorithmic description of the infinite set of natural numbers. Spaces produced by the small object argument are infamously complicated and difficult to understand; the type-theoretic approach is potentially much simpler, bypassing the need for any explicit construction

by giving direct access to the appropriate “induction principle”. Thus, the combination of univalence and higher inductive types suggests the possibility of a revolution, of sorts, in the practice of homotopy theory.

Sets in univalent foundations

We have claimed that univalent foundations can eventually serve as a foundation for “all” of mathematics, but so far we have discussed only homotopy theory. Of course, there are many specific examples of the use of type theory without the new homotopy type theory features to formalize mathematics, such as the recent formalization of the Feit–Thompson odd-order theorem in Coq [GAA⁺13].

But the traditional view is that mathematics is founded on set theory, in the sense that all mathematical objects and constructions can be coded into a theory such as Zermelo–Fraenkel set theory (ZF). However, it is well-established by now that for most mathematics outside of set theory proper, the intricate hierarchical membership structure of sets in ZF is really unnecessary: a more “structural” theory, such as Lawvere’s Elementary Theory of the Category of Sets [Law05], suffices.

In univalent foundations, the basic objects are “homotopy types” rather than sets, but we can *define* a class of types which behave like sets. Homotopically, these can be thought of as spaces in which every connected component is contractible, i.e. those which are homotopy equivalent to a discrete space. It is a theorem that the category of such “sets” satisfies Lawvere’s axioms (or related ones, depending on the details of the theory). Thus, any sort of mathematics that can be represented in an ETCS-like theory (which, experience suggests, is essentially all of mathematics) can equally well be represented in univalent foundations.

This supports the claim that univalent foundations is at least as good as existing foundations of mathematics. A mathematician working in univalent foundations can build structures out of sets in a familiar way, with more general homotopy types waiting in the foundational background until there is need of them. For this reason, most of the applications in this book have been chosen to be areas where univalent foundations has something *new* to contribute that distinguishes it from existing foundational systems.

Unsurprisingly, homotopy theory and category theory are two of these, but perhaps less obvious is that univalent foundations has something new and interesting to offer even in subjects such as set theory and real analysis. For instance, the univalence axiom allows us to identify isomorphic structures, while higher inductive types allow direct descriptions of objects by their universal properties. Thus we can generally avoid resorting to arbitrarily chosen representatives or transfinite iterative constructions. In fact, even the objects of study in ZF set theory can be characterized, inside the sets of univalent foundations, by such an inductive universal property.

Informal type theory

One difficulty often encountered by the classical mathematician when faced with learning about type theory is that it is usually presented as a fully or partially formalized deductive system. This style, which is very useful for proof-theoretic investigations, is not particularly convenient for use in applied, informal reasoning. Nor is it even familiar to most working mathematicians, even those who might be interested in foundations of mathematics. One objective of the present work is to develop an informal style of doing mathematics in univalent foundations that is at once rigorous and precise, but is also closer to the language and style of presentation of everyday

mathematics.

In present-day mathematics, one usually constructs and reasons about mathematical objects in a way that could in principle, one presumes, be formalized in a system of elementary set theory, such as ZFC — at least given enough ingenuity and patience. For the most part, one does not even need to be aware of this possibility, since it largely coincides with the condition that a proof be “fully rigorous” (in the sense that all mathematicians have come to understand intuitively through education and experience). But one does need to learn to be careful about a few aspects of “informal set theory”: the use of collections too large or inchoate to be sets; the axiom of choice and its equivalents; even (for undergraduates) the method of proof by contradiction; and so on. Adopting a new foundational system such as homotopy type theory as the *implicit formal basis* of informal reasoning will require adjusting some of one’s instincts and practices. The present text is intended to serve as an example of this “new kind of mathematics”, which is still informal, but could now in principle be formalized in homotopy type theory, rather than ZFC, again given enough ingenuity and patience.

It is worth emphasizing that, in this new system, such formalization can have real practical benefits. The formal system of type theory is suited to computer systems and has been implemented in existing proof assistants. A proof assistant is a computer program which guides the user in construction of a fully formal proof, only allowing valid steps of reasoning. It also provides some degree of automation, can search libraries for existing theorems, and can even extract numerical algorithms from the resulting (constructive) proofs.

We believe that this aspect of the univalent foundations program distinguishes it from other approaches to foundations, potentially providing a new practical utility for the working mathematician. Indeed, proof assistants based on older type theories have already been used to formalize substantial mathematical proofs, such as the four-color theorem and the Feit–Thompson theorem. Computer implementations of univalent foundations are presently works in progress (like the theory itself). However, even its currently available implementations (which are mostly small modifications to existing proof assistants such as COQ and AGDA) have already demonstrated their worth, not only in the formalization of known proofs, but in the discovery of new ones. Indeed, many of the proofs described in this book were actually *first* done in a fully formalized form in a proof assistant, and are only now being “unformalized” for the first time — a reversal of the usual relation between formal and informal mathematics.

One can imagine a not-too-distant future when it will be possible for mathematicians to verify the correctness of their own papers by working within the system of univalent foundations, formalized in a proof assistant, and that doing so will become as natural as typesetting their own papers in \TeX . In principle, this could be equally true for any other foundational system, but we believe it to be more practically attainable using univalent foundations, as witnessed by the present work and its formal counterpart.

Constructivity

One of the most striking differences between classical foundations and type theory is the idea of *proof relevance*, according to which mathematical statements, and even their proofs, become first-class mathematical objects. In type theory, we represent mathematical statements by types, which can be regarded simultaneously as both mathematical constructions and mathematical assertions, a conception also known as *propositions as types*. Accordingly, we can regard a term $a : A$ as both an element of the type A (or in homotopy type theory, a point of the space A), and at the same time, a proof of the proposition A . To take an example, suppose we have sets A and

B (discrete spaces), and consider the statement “ A is isomorphic to B ”. In type theory, this can be rendered as:

$$\text{Iso}(A, B) := \sum_{(f:A \rightarrow B)} \sum_{(g:B \rightarrow A)} \left((\Pi_{(x:A)} g(f(x)) = x) \times (\Pi_{(y:B)} f(g(y)) = y) \right).$$

Reading the type constructors Σ, Π, \times here as “there exists”, “for all”, and “and” respectively yields the usual formulation of “ A and B are isomorphic”; on the other hand, reading them as sums and products yields the *type of all isomorphisms* between A and B ! To prove that A and B are isomorphic, one constructs a proof $p : \text{Iso}(A, B)$, which is therefore the same as constructing an isomorphism between A and B , i.e., exhibiting a pair of functions f, g together with *proofs* that their composites are the respective identity maps. The latter proofs, in turn, are nothing but homotopies of the appropriate sorts. In this way, *proving a proposition is the same as constructing an element of some particular type*. In particular, to prove a statement of the form “ A and B ” is just to prove A and to prove B , i.e., to give an element of the type $A \times B$. And to prove that A implies B is just to find an element of $A \rightarrow B$, i.e. a function from A to B (determining a mapping of proofs of A to proofs of B).

The logic of propositions-as-types is flexible and supports many variations, such as using only a subclass of types to represent propositions. In homotopy type theory, there are natural such subclasses arising from the fact that the system of all types, just like spaces in classical homotopy theory, is “stratified” according to the dimensions in which their higher homotopy structure exists or collapses. In particular, Voevodsky has found a purely type-theoretic definition of *homotopy n-types*, corresponding to spaces with no nontrivial homotopy information above dimension n . (The 0-types are the “sets” mentioned previously as satisfying Lawvere’s axioms.) Moreover, with higher inductive types, we can universally “truncate” a type into an n -type; in classical homotopy theory this would be its n^{th} Postnikov section. Particularly important for logic is the case of homotopy (-1) -types, which we call *mere propositions*. Classically, every (-1) -type is empty or contractible; we interpret these possibilities as the truth values “false” and “true” respectively.

Using all types as propositions yields a very “constructive” conception of logic; for more on this, see [Kol32, TvD88a, TvD88b]. For instance, every proof that something exists carries with it enough information to actually find such an object; and every proof that “ A or B ” holds is either a proof that A holds or a proof that B holds. Thus, from every proof we can automatically extract an algorithm; this can be very useful in applications to computer programming.

On the other hand, however, this logic does diverge from the traditional understanding of existence proofs in mathematics. In particular, it does not faithfully represent certain important classical principles of reasoning, such as the axiom of choice and the law of excluded middle. For these we need to use the “ (-1) -truncated” logic, in which only the homotopy (-1) -types represent propositions.

More specifically, consider on one hand the *axiom of choice*: “if for every $x : A$ there exists a $y : B$ such that $R(x, y)$, there is a function $f : A \rightarrow B$ such that for all $x : A$ we have $R(x, f(x))$.” The pure propositions-as-types notion of “there exists” is strong enough to make this statement simply provable — yet it does not have all the consequences of the usual axiom of choice. However, in (-1) -truncated logic, this statement is not automatically true, but is a strong assumption with the same sorts of consequences as its counterpart in classical set theory.

On the other hand, consider the *law of excluded middle*: “for all A , either A or not A .” Interpreting this in the pure propositions-as-types logic yields a statement that is inconsistent with the univalence axiom. For since proving “ A ” means exhibiting an element of it, this assumption

would give a uniform way of selecting an element from every nonempty type — a sort of Hilbertian choice operator. Univalence implies that the element of A selected by such a choice operator must be invariant under all self-equivalences of A , since these are identified with self-identities and every operation must respect identity; but clearly some types have automorphisms with no fixed points, e.g. we can swap the elements of a two-element type. However, the “ (-1) -truncated law of excluded middle”, though also not automatically true, may consistently be assumed with most of the same consequences as in classical mathematics.

In other words, while the pure propositions-as-types logic is “constructive” in the strong algorithmic sense mentioned above, the default (-1) -truncated logic is “constructive” in a different sense (namely, that of the logic formalized by Heyting under the name “intuitionistic”); and to the latter we may freely add the axioms of choice and excluded middle to obtain a logic that may be called “classical”. Thus, homotopy type theory is compatible with both constructive and classical conceptions of logic, and many more besides. Indeed, the homotopical perspective reveals that classical and constructive logic can coexist, as endpoints of a spectrum of different systems, with an infinite number of possibilities in between (the homotopy n -types for $-1 < n < \infty$). We may speak of “ LEM_n ” and “ AC_n ”, with AC_∞ being provable and LEM_∞ inconsistent with univalence, while AC_{-1} and LEM_{-1} are the versions familiar to classical mathematicians (hence in most cases it is appropriate to assume the subscript (-1) when none is given). Indeed, one can even have useful systems in which only *certain* types satisfy such further “classical” principles, while types in general remain “constructive”.

It is worth emphasizing that univalent foundations does not *require* the use of constructive or intuitionistic logic. Most of classical mathematics which depends on the law of excluded middle and the axiom of choice can be performed in univalent foundations, simply by assuming that these two principles hold (in their proper, (-1) -truncated, form). However, type theory does encourage avoiding these principles when they are unnecessary, for several reasons.

First of all, every mathematician knows that a theorem is more powerful when proven using fewer assumptions, since it applies to more examples. The situation with AC and LEM is no different: type theory admits many interesting “nonstandard” models, such as in sheaf toposes, where classicality principles such as AC and LEM tend to fail. Homotopy type theory admits similar models in higher toposes, such as are studied in [TV02, Rez05, Lur09]. Thus, if we avoid using these principles, the theorems we prove will be valid internally to all such models.

Secondly, one of the additional virtues of type theory is its computable character. In addition to being a foundation for mathematics, type theory is a formal theory of computation, and can be treated as a powerful programming language. From this perspective, the rules of the system cannot be chosen arbitrarily the way set-theoretic axioms can: there must be a harmony between them which allows all proofs to be “executed” as programs. We do not yet fully understand the new principles introduced by homotopy type theory, such as univalence and higher inductive types, from this point of view, but the basic outlines are emerging; see, for example, [LH12]. It has been known for a long time, however, that principles such as AC and LEM are fundamentally antithetical to computability, since they assert baldly that certain things exist without giving any way to compute them. Thus, avoiding them is necessary to maintain the character of type theory as a theory of computation.

Fortunately, constructive reasoning is not as hard as it may seem. In some cases, simply by rephrasing some definitions, a theorem can be made constructive and its proof more elegant. Moreover, in univalent foundations this seems to happen more often. For instance:

- (i) In set-theoretic foundations, at various points in homotopy theory and category theory one

needs the axiom of choice to perform transfinite constructions. But with higher inductive types, we can encode these constructions directly and constructively. In particular, none of the “synthetic” homotopy theory in Chapter 8 requires LEM or AC.

- (ii) In set-theoretic foundations, the statement “every fully faithful and essentially surjective functor is an equivalence of categories” is equivalent to the axiom of choice. But with the univalence axiom, it is just *true*; see Chapter 9.
- (iii) In set theory, various circumlocutions are required to obtain notions of “cardinal number” and “ordinal number” which canonically represent isomorphism classes of sets and well-ordered sets, respectively — possibly involving the axiom of choice or the axiom of foundation. But with univalence and higher inductive types, we can obtain such representatives directly by truncating the universe; see Chapter 10.
- (iv) In set-theoretic foundations, the definition of the real numbers as equivalence classes of Cauchy sequences requires either the law of excluded middle or the axiom of (countable) choice to be well-behaved. But with higher inductive types, we can give a version of this definition which is well-behaved and avoids any choice principles; see Chapter 11.

Of course, these simplifications could as well be taken as evidence that the new methods will not, ultimately, prove to be really constructive. However, we emphasize again that the reader does not have to care, or worry, about constructivity in order to read this book. The point is that in all of the above examples, the version of the theory we give has independent advantages, whether or not LEM and AC are assumed to be available. Constructivity, if attained, will be an added bonus.

Given this discussion of adding new principles such as univalence, higher inductive types, AC, and LEM, one may wonder whether the resulting system remains consistent. (One of the original virtues of type theory, relative to set theory, was that it can be seen to be consistent by proof-theoretic means). As with any foundational system, consistency is a relative question: “consistent with respect to what?” The short answer is that all of the constructions and axioms considered in this book have a model in the category of Kan complexes, due to Voevodsky [KLV12] (see [LS17] for higher inductive types). Thus, they are known to be consistent relative to ZFC (with as many inaccessible cardinals as we need nested univalent universes). Giving a more traditionally type-theoretic account of this consistency is work in progress (see, e.g., [LH12, BCH13]).

We summarize the different points of view of the type-theoretic operations in Table 1.

Open problems

For those interested in contributing to this new branch of mathematics, it may be encouraging to know that there are many interesting open questions.

Perhaps the most pressing of them is the “constructivity” of the Univalence Axiom, posed by Voevodsky in [Voe12]. The basic system of type theory follows the structure of Gentzen’s natural deduction. Logical connectives are defined by their introduction rules, and have elimination rules justified by computation rules. Following this pattern, and using Tait’s computability method, originally designed to analyse Gödel’s *Dialectica* interpretation, one can show the property of *normalization* for type theory. This in turn implies important properties such as decidability of type-checking (a crucial property since type-checking corresponds to proof-checking, and one can argue that we should be able to “recognize a proof when we see one”), and the so-called “canonicity property” that any closed term of the type of natural numbers reduces to a

Types	Logic	Sets	Homotopy
A	proposition	set	space
$a : A$	proof	element	point
$B(x)$	predicate	family of sets	fibration
$b(x) : B(x)$	conditional proof	family of elements	section
$\mathbf{0}, \mathbf{1}$	\perp, \top	$\emptyset, \{\emptyset\}$	$\emptyset, *$
$A + B$	$A \vee B$	disjoint union	coproduct
$A \times B$	$A \wedge B$	set of pairs	product space
$A \rightarrow B$	$A \Rightarrow B$	set of functions	function space
$\sum_{(x:A)} B(x)$	$\exists_{x:A} B(x)$	disjoint sum	total space
$\prod_{(x:A)} B(x)$	$\forall_{x:A} B(x)$	product	space of sections
id_A	equality =	$\{ (x, x) \mid x \in A \}$	path space A^I

Table 1: Comparing points of view on type-theoretic operations

numeral. This last property, and the uniform structure of introduction/elimination rules, are lost when one extends type theory with an axiom, such as the axiom of function extensionality, or the univalence axiom. Voevodsky has formulated a precise mathematical conjecture connected to this question of canonicity for type theory extended with the axiom of Univalence: given a closed term of the type of natural numbers, is it always possible to find a numeral and a proof that this term is equal to this numeral, where this proof of equality may itself use the univalence axiom? More generally, an important issue is whether it is possible to provide a constructive justification of the univalence axiom. What about if one adds other homotopically motivated constructions, like higher inductive types? These questions remain open at the present time, although methods are currently being developed to try to find answers.

Another basic issue is the difficulty of working with types, such as the natural numbers, that are essentially sets (i.e., discrete spaces), containing only trivial paths. At present, homotopy type theory can really only characterize spaces up to homotopy equivalence, which means that these “discrete spaces” may only be *homotopy equivalent* to discrete spaces. Type-theoretically, this means there are many paths that are equal to reflexivity, but not *judgmentally* equal to it (see §1.1 for the meaning of “judgmentally”). While this homotopy-invariance has advantages, these “meaningless” identity terms do introduce needless complications into arguments and constructions, so it would be convenient to have a systematic way of eliminating or collapsing them.

A more specialized, but no less important, problem is the relation between homotopy type theory and the research on *higher toposes* currently happening at the intersection of higher category theory and homotopy theory. There is a growing conviction among those familiar with both subjects that they are intimately connected. For instance, the notion of a univalent universe should coincide with that of an object classifier, while higher inductive types should be an “elementary” reflection of local presentability. More generally, homotopy type theory should be the “internal language” of $(\infty, 1)$ -toposes, just as intuitionistic higher-order logic is the internal language of ordinary 1-toposes. Despite this general consensus, however, details remain to be worked out — in particular, questions of coherence and strictness remain to be addressed — and doing so will undoubtedly lead to further insights into both concepts.

But by far the largest field of work to be done is in the ongoing formalization of everyday mathematics in this new system. Recent successes in formalizing some facts from basic homotopy theory and category theory have been encouraging; some of these are described in Chapters 8 and 9. Obviously, however, much work remains to be done.

The homotopy type theory community maintains a web site and group blog at <http://homotopytypetheory.org>, as well as a discussion email list. Newcomers are always welcome!

How to read this book

This book is divided into two parts. Part I, “Foundations”, develops the fundamental concepts of homotopy type theory. This is the mathematical foundation on which the development of specific subjects is built, and which is required for the understanding of the univalent foundations approach. To a programmer, this is “library code”. Since univalent foundations is a new and different kind of mathematics, its basic notions take some getting used to; thus Part I is fairly extensive.

Part II, “Mathematics”, consists of four chapters that build on the basic notions of Part I to exhibit some of the new things we can do with univalent foundations in four different areas of mathematics: homotopy theory (Chapter 8), category theory (Chapter 9), set theory (Chapter 10), and real analysis (Chapter 11). The chapters in Part II are more or less independent of each other, although occasionally one will use a lemma proven in another.

A reader who wants to seriously understand univalent foundations, and be able to work in it, will eventually have to read and understand most of Part I. However, a reader who just wants to get a taste of univalent foundations and what it can do may understandably balk at having to work through over 200 pages before getting to the “meat” in Part II. Fortunately, not all of Part I is necessary in order to read the chapters in Part II. Each chapter in Part II begins with a brief overview of its subject, what univalent foundations has to contribute to it, and the necessary background from Part I, so the courageous reader can turn immediately to the appropriate chapter for their favorite subject. For those who want to understand one or more chapters in Part II more deeply than this, but are not ready to read all of Part I, we provide here a brief summary of Part I, with remarks about which parts are necessary for which chapters in Part II.

Chapter 1 is about the basic notions of type theory, prior to any homotopical interpretation. A reader who is familiar with Martin-Löf type theory can quickly skim it to pick up the particulars of the theory we are using. However, readers without experience in type theory will need to read Chapter 1, as there are many subtle differences between type theory and other foundations such as set theory.

Chapter 2 introduces the homotopical viewpoint on type theory, along with the basic notions supporting this view, and describes the homotopical behavior of each component of the type theory from Chapter 1. It also introduces the *univalence axiom* (§2.10) — the first of the two basic innovations of homotopy type theory. Thus, it is quite basic and we encourage everyone to read it, especially §§2.1–2.4.

Chapter 3 describes how we represent logic in homotopy type theory, and its connection to classical logic as well as to constructive and intuitionistic logic. Here we define the law of excluded middle, the axiom of choice, and the axiom of propositional resizing (although, for the most part, we do not need to assume any of these in the rest of the book), as well as the *propositional truncation* which is essential for representing traditional logic. This chapter is essential background for Chapters 10 and 11, less important for Chapter 9, and not so necessary for

Chapter 8.

Chapters 4 and 5 study two special topics in detail: equivalences (and related notions) and generalized inductive definitions. While these are important subjects in their own rights and provide a deeper understanding of homotopy type theory, for the most part they are not necessary for Part II. Only a few lemmas from Chapter 4 are used here and there, while the general discussions in §§5.1, 5.6 and 5.7 are helpful for providing the intuition required for Chapter 6. The generalized sorts of inductive definition discussed in §5.7 are also used in a few places in Chapters 10 and 11.

Chapter 6 introduces the second basic innovation of homotopy type theory — *higher inductive types* — with many examples. Higher inductive types are the primary object of study in Chapter 8, and some particular ones play important roles in Chapters 10 and 11. They are not so necessary for Chapter 9, although one example is used in §9.9.

Finally, Chapter 7 discusses homotopy n -types and related notions such as n -connected types. These notions are important for Chapter 8, but not so important in the rest of Part II, although the case $n = -1$ of some of the lemmas are used in §10.1.

This completes Part I. As mentioned above, Part II consists of four largely unrelated chapters, each describing what univalent foundations has to offer to a particular subject.

Of the chapters in Part II, Chapter 8 (Homotopy theory) is perhaps the most radical. Univalent foundations has a very different “synthetic” approach to homotopy theory in which homotopy types are the basic objects (namely, the types) rather than being constructed using topological spaces or some other set-theoretic model. This enables new styles of proof for classical theorems in algebraic topology, of which we present a sampling, from $\pi_1(\mathbb{S}^1) = \mathbb{Z}$ to the Freudenthal suspension theorem.

In Chapter 9 (Category theory), we develop some basic (1-)category theory, adhering to the principle of the univalence axiom that *equality is isomorphism*. This has the pleasant effect of ensuring that all definitions and constructions are automatically invariant under equivalence of categories: indeed, equivalent categories are equal just as equivalent types are equal. (It also has connections to higher category theory and higher topos theory.)

Chapter 10 (Set theory) studies sets in univalent foundations. The category of sets has its usual properties, hence provides a foundation for any mathematics that doesn’t need homotopical or higher-categorical structures. We also observe that univalence makes cardinal and ordinal numbers a bit more pleasant, and that higher inductive types yield a cumulative hierarchy satisfying the usual axioms of Zermelo–Fraenkel set theory.

In Chapter 11 (Real numbers), we summarize the construction of Dedekind real numbers, and then observe that higher inductive types allow a definition of Cauchy real numbers that avoids some associated problems in constructive mathematics. Then we sketch a similar approach to Conway’s surreal numbers.

Each chapter in this book ends with a Notes section, which collects historical comments, references to the literature, and attributions of results, to the extent possible. We have also included Exercises at the end of each chapter, to assist the reader in gaining familiarity with doing mathematics in univalent foundations.

Finally, recall that this book was written as a massively collaborative effort by a large number of people. We have done our best to achieve consistency in terminology and notation, and to put the mathematics in a linear sequence that flows logically, but it is very likely that some imperfections remain. We ask the reader’s forgiveness for any such infelicities, and welcome suggestions for improvement of the next edition.

PART I

FOUNDATIONS

Chapter 1

Type theory

1.1 Type theory versus set theory

Homotopy type theory is (among other things) a foundational language for mathematics, i.e., an alternative to Zermelo–Fraenkel set theory. However, it behaves differently from set theory in several important ways, and that can take some getting used to. Explaining these differences carefully requires us to be more formal here than we will be in the rest of the book. As stated in the introduction, our goal is to write type theory *informally*; but for a mathematician accustomed to set theory, more precision at the beginning can help avoid some common misconceptions and mistakes.

We note that a set-theoretic foundation has two “layers”: the deductive system of first-order logic, and, formulated inside this system, the axioms of a particular theory, such as ZFC. Thus, set theory is not only about sets, but rather about the interplay between sets (the objects of the second layer) and propositions (the objects of the first layer).

By contrast, type theory is its own deductive system: it need not be formulated inside any superstructure, such as first-order logic. Instead of the two basic notions of set theory, sets and propositions, type theory has one basic notion: *types*. Propositions (statements which we can prove, disprove, assume, negate, and so on¹) are identified with particular types, via the correspondence shown in Table 1 on page 11. Thus, the mathematical activity of *proving a theorem* is identified with a special case of the mathematical activity of *constructing an object*—in this case, an inhabitant of a type that represents a proposition.

This leads us to another difference between type theory and set theory, but to explain it we must say a little about deductive systems in general. Informally, a deductive system is a collection of **rules** for deriving things called **judgments**. If we think of a deductive system as a formal game, then the judgments are the “positions” in the game which we reach by following the game rules. We can also think of a deductive system as a sort of algebraic theory, in which case the judgments are the elements (like the elements of a group) and the deductive rules are the operations (like the group multiplication). From a logical point of view, the judgments can be considered to be the “external” statements, living in the metatheory, as opposed to the “internal” statements of the theory itself.

¹Confusingly, it is also a common practice (dating back to Euclid) to use the word “proposition” synonymously with “theorem”. We will confine ourselves to the logician’s usage, according to which a *proposition* is a statement *susceptible to proof*, whereas a *theorem* (or “lemma” or “corollary”) is such a statement that *has been proven*. Thus “ $0 = 1$ ” and its negation “ $\neg(0 = 1)$ ” are both propositions, but only the latter is a theorem.

In the deductive system of first-order logic (on which set theory is based), there is only one kind of judgment: that a given proposition has a proof. That is, each proposition A gives rise to a judgment “ A has a proof”, and all judgments are of this form. A rule of first-order logic such as “from A and B infer $A \wedge B$ ” is actually a rule of “proof construction” which says that given the judgments “ A has a proof” and “ B has a proof”, we may deduce that “ $A \wedge B$ has a proof”. Note that the judgment “ A has a proof” exists at a different level from the *proposition* A itself, which is an internal statement of the theory.

The basic judgment of type theory, analogous to “ A has a proof”, is written “ $a : A$ ” and pronounced as “the term a has type A ”, or more loosely “ a is an element of A ” (or, in homotopy type theory, “ a is a point of A ”). When A is a type representing a proposition, then a may be called a *witness* to the provability of A , or *evidence* of the truth of A (or even a *proof* of A , but we will try to avoid this confusing terminology). In this case, the judgment $a : A$ is derivable in type theory (for some a) precisely when the analogous judgment “ A has a proof” is derivable in first-order logic (modulo differences in the axioms assumed and in the encoding of mathematics, as we will discuss throughout the book).

On the other hand, if the type A is being treated more like a set than like a proposition (although as we will see, the distinction can become blurry), then “ $a : A$ ” may be regarded as analogous to the set-theoretic statement “ $a \in A$ ”. However, there is an essential difference in that “ $a : A$ ” is a *judgment* whereas “ $a \in A$ ” is a *proposition*. In particular, when working internally in type theory, we cannot make statements such as “if $a : A$ then it is not the case that $b : B$ ”, nor can we “disprove” the judgment “ $a : A$ ”.

A good way to think about this is that in set theory, “membership” is a relation which may or may not hold between two pre-existing objects “ a ” and “ A ”, while in type theory we cannot talk about an element “ a ” in isolation: every element *by its very nature* is an element of some type, and that type is (generally speaking) uniquely determined. Thus, when we say informally “let x be a natural number”, in set theory this is shorthand for “let x be a thing and assume that $x \in \mathbb{N}$ ”, whereas in type theory “let $x : \mathbb{N}$ ” is an atomic statement: we cannot introduce a variable without specifying its type.

At first glance, this may seem an uncomfortable restriction, but it is arguably closer to the intuitive mathematical meaning of “let x be a natural number”. In practice, it seems that whenever we actually *need* “ $a \in A$ ” to be a proposition rather than a judgment, there is always an ambient set B of which a is known to be an element and A is known to be a subset. This situation is also easy to represent in type theory, by taking a to be an element of the type B , and A to be a predicate on B ; see §3.5.

A last difference between type theory and set theory is the treatment of equality. The familiar notion of equality in mathematics is a proposition: e.g. we can disprove an equality or assume an equality as a hypothesis. Since in type theory, propositions are types, this means that equality is a type: for elements $a, b : A$ (that is, both $a : A$ and $b : A$) we have a type “ $a =_A b$ ”. (In *homotopy* type theory, of course, this equality proposition can behave in unfamiliar ways: see §1.12 and Chapter 2, and the rest of the book). When $a =_A b$ is inhabited, we say that a and b are **(propositionally) equal**.

However, in type theory there is also a need for an equality *judgment*, existing at the same level as the judgment “ $x : A$ ”. This is called **judgmental equality** or **definitional equality**, and we write it as $a \equiv b : A$ or simply $a \equiv b$. It is helpful to think of this as meaning “equal by definition”. For instance, if we define a function $f : \mathbb{N} \rightarrow \mathbb{N}$ by the equation $f(x) = x^2$, then the expression $f(3)$ is equal to 3^2 *by definition*. Inside the theory, it does not make sense to negate or assume an equality-by-definition; we cannot say “if x is equal to y by definition, then z is

not equal to w by definition". Whether or not two expressions are equal by definition is just a matter of expanding out the definitions; in particular, it is algorithmically decidable (though the algorithm is necessarily meta-theoretic, not internal to the theory).

As type theory becomes more complicated, judgmental equality can get more subtle than this, but it is a good intuition to start from. Alternatively, if we regard a deductive system as an algebraic theory, then judgmental equality is simply the equality in that theory, analogous to the equality between elements of a group—the only potential for confusion is that there is *also* an object *inside* the deductive system of type theory (namely the type " $a = b$ ") which behaves internally as a notion of "equality".

The reason we *want* a judgmental notion of equality is so that it can control the other form of judgment, " $a : A$ ". For instance, suppose we have given a proof that $3^2 = 9$, i.e. we have derived the judgment $p : (3^2 = 9)$ for some p . Then the same witness p ought to count as a proof that $f(3) = 9$, since $f(3)$ is 3^2 *by definition*. The best way to represent this is with a rule saying that given the judgments $a : A$ and $A \equiv B$, we may derive the judgment $a : B$.

Thus, for us, type theory will be a deductive system based on two forms of judgment:

Judgment	Meaning
$a : A$	" a is an object of type A "
$a \equiv b : A$	" a and b are definitionally equal objects of type A "

When introducing a definitional equality, i.e., defining one thing to be equal to another, we will use the symbol " \equiv ". Thus, the above definition of the function f would be written as $f(x) \equiv x^2$.

Because judgments cannot be put together into more complicated statements, the symbols ":" and " \equiv " bind more loosely than anything else.² Thus, for instance, " $p : x = y$ " should be parsed as " $p : (x = y)$ ", which makes sense since " $x = y$ " is a type, and not as " $(p : x) = y$ ", which is senseless since " $p : x$ " is a judgment and cannot be equal to anything. Similarly, " $A \equiv x = y$ " can only be parsed as " $A \equiv (x = y)$ ", although in extreme cases such as this, one ought to add parentheses anyway to aid reading comprehension. Moreover, later on we will fall into the common notation of chaining together equalities — e.g. writing $a = b = c = d$ to mean " $a = b$ and $b = c$ and $c = d$, hence $a = d$ " — and we will also include judgmental equalities in such chains. Context usually suffices to make the intent clear.

This is perhaps also an appropriate place to mention that the common mathematical notation " $f : A \rightarrow B$ ", expressing the fact that f is a function from A to B , can be regarded as a typing judgment, since we use " $A \rightarrow B$ " as notation for the type of functions from A to B (as is standard practice in type theory; see §1.4).

Judgments may depend on *assumptions* of the form $x : A$, where x is a variable and A is a type. For example, we may construct an object $m + n : \mathbb{N}$ under the assumptions that $m, n : \mathbb{N}$. Another example is that assuming A is a type, $x, y : A$, and $p : x =_A y$, we may construct an element $p^{-1} : y =_A x$. The collection of all such assumptions is called the **context**; from a topological point of view it may be thought of as a "parameter space". In fact, technically the context must be an ordered list of assumptions, since later assumptions may depend on previous ones: the assumption $x : A$ can only be made *after* the assumptions of any variables appearing in the type A .

²In formalized type theory, commas and turnstiles can bind even more loosely. For instance, $x : A, y : B \vdash c : C$ is parsed as $((x : A), (y : B)) \vdash (c : C)$. However, in this book we refrain from such notation until ??.

If the type A in an assumption $x : A$ represents a proposition, then the assumption is a type-theoretic version of a *hypothesis*: we assume that the proposition A holds. When types are regarded as propositions, we may omit the names of their proofs. Thus, in the second example above we may instead say that assuming $x =_A y$, we can prove $y =_A x$. However, since we are doing “proof-relevant” mathematics, we will frequently refer back to proofs as objects. In the example above, for instance, we may want to establish that p^{-1} together with the proofs of transitivity and reflexivity behave like a groupoid; see Chapter 2.

Note that under this meaning of the word *assumption*, we can assume a propositional equality (by assuming a variable $p : x = y$), but we cannot assume a judgmental equality $x \equiv y$, since it is not a type that can have an element. However, we can do something else which looks kind of like assuming a judgmental equality: if we have a type or an element which involves a variable $x : A$, then we can *substitute* any particular element $a : A$ for x to obtain a more specific type or element. We will sometimes use language like “now assume $x \equiv a$ ” to refer to this process of substitution, even though it is not an *assumption* in the technical sense introduced above.

By the same token, we cannot *prove* a judgmental equality either, since it is not a type in which we can exhibit a witness. Nevertheless, we will sometimes state judgmental equalities as part of a theorem, e.g. “there exists $f : A \rightarrow B$ such that $f(x) \equiv y$ ”. This should be regarded as the making of two separate judgments: first we make the judgment $f : A \rightarrow B$ for some element f , then we make the additional judgment that $f(x) \equiv y$.

In the rest of this chapter, we attempt to give an informal presentation of type theory, sufficient for the purposes of this book; we give a more formal account in ???. Aside from some fairly obvious rules (such as the fact that judgmentally equal things can always be substituted for each other), the rules of type theory can be grouped into *type formers*. Each type former consists of a way to construct types (possibly making use of previously constructed types), together with rules for the construction and behavior of elements of that type. In most cases, these rules follow a fairly predictable pattern, but we will not attempt to make this precise here; see however the beginning of §1.5 and also Chapter 5.

An important aspect of the type theory presented in this chapter is that it consists entirely of *rules*, without any *axioms*. In the description of deductive systems in terms of judgments, the *rules* are what allow us to conclude one judgment from a collection of others, while the *axioms* are the judgments we are given at the outset. If we think of a deductive system as a formal game, then the rules are the rules of the game, while the axioms are the starting position. And if we think of a deductive system as an algebraic theory, then the rules are the operations of the theory, while the axioms are the *generators* for some particular free model of that theory.

In set theory, the only rules are the rules of first-order logic (such as the rule allowing us to deduce “ $A \wedge B$ has a proof” from “ A has a proof” and “ B has a proof”): all the information about the behavior of sets is contained in the axioms. By contrast, in type theory, it is usually the *rules* which contain all the information, with no axioms being necessary. For instance, in §1.5 we will see that there is a rule allowing us to deduce the judgment “ $(a, b) : A \times B$ ” from “ $a : A$ ” and “ $b : B$ ”, whereas in set theory the analogous statement would be (a consequence of) the pairing axiom.

The advantage of formulating type theory using only rules is that rules are “procedural”. In particular, this property is what makes possible (though it does not automatically ensure) the good computational properties of type theory, such as “canonicity”. However, while this style works for traditional type theories, we do not yet understand how to formulate everything we need for *homotopy* type theory in this way. In particular, in §§2.9 and 2.10 and Chapter 6 we will have to augment the rules of type theory presented in this chapter by introducing additional ax-

ioms, notably the *univalence axiom*. In this chapter, however, we confine ourselves to a traditional rule-based type theory.

1.2 Function types

Given types A and B , we can construct the type $A \rightarrow B$ of **functions** with domain A and codomain B . We also sometimes refer to functions as **maps**. Unlike in set theory, functions are not defined as functional relations; rather they are a primitive concept in type theory. We explain the function type by prescribing what we can do with functions, how to construct them and what equalities they induce.

Given a function $f : A \rightarrow B$ and an element of the domain $a : A$, we can **apply** the function to obtain an element of the codomain B , denoted $f(a)$ and called the **value** of f at a . It is common in type theory to omit the parentheses and denote $f(a)$ simply by $f a$, and we will sometimes do this as well.

But how can we construct elements of $A \rightarrow B$? There are two equivalent ways: either by direct definition or by using λ -abstraction. Introducing a function by definition means that we introduce a function by giving it a name — let's say, f — and saying we define $f : A \rightarrow B$ by giving an equation

$$f(x) : \equiv \Phi \tag{1.2.1}$$

where x is a variable and Φ is an expression which may use x . In order for this to be valid, we have to check that $\Phi : B$ assuming $x : A$.

Now we can compute $f(a)$ by replacing the variable x in Φ with a . As an example, consider the function $f : \mathbb{N} \rightarrow \mathbb{N}$ which is defined by $f(x) : \equiv x + x$. (We will define \mathbb{N} and $+$ in §1.9.) Then $f(2)$ is judgmentally equal to $2 + 2$.

If we don't want to introduce a name for the function, we can use **λ -abstraction**. Given an expression Φ of type B which may use $x : A$, as above, we write $\lambda(x : A). \Phi$ to indicate the same function defined by (1.2.1). Thus, we have

$$(\lambda(x : A). \Phi) : A \rightarrow B.$$

For the example in the previous paragraph, we have the typing judgment

$$(\lambda(x : \mathbb{N}). x + x) : \mathbb{N} \rightarrow \mathbb{N}.$$

As another example, for any types A and B and any element $y : B$, we have a **constant function** $(\lambda(x : A). y) : A \rightarrow B$.

We generally omit the type of the variable x in a λ -abstraction and write $\lambda x. \Phi$, since the typing $x : A$ is inferable from the judgment that the function $\lambda x. \Phi$ has type $A \rightarrow B$. By convention, the “scope” of the variable binding “ $\lambda x.$ ” is the entire rest of the expression, unless delimited with parentheses. Thus, for instance, $\lambda x. x + x$ should be parsed as $\lambda x. (x + x)$, not as $(\lambda x. x) + x$ (which would, in this case, be ill-typed anyway).

Another equivalent notation is

$$(x \mapsto \Phi) : A \rightarrow B.$$

We may also sometimes use a blank “ $-$ ” in the expression Φ in place of a variable, to denote an implicit λ -abstraction. For instance, $g(x, -)$ is another way to write $\lambda y. g(x, y)$.

Now a λ -abstraction is a function, so we can apply it to an argument $a : A$. We then have the following **computation rule**³, which is a definitional equality:

$$(\lambda x. \Phi)(a) \equiv \Phi'$$

where Φ' is the expression Φ in which all occurrences of x have been replaced by a . Continuing the above example, we have

$$(\lambda x. x + x)(2) \equiv 2 + 2.$$

Note that from any function $f : A \rightarrow B$, we can construct a lambda abstraction function $\lambda x. f(x)$. Since this is by definition “the function that applies f to its argument” we consider it to be definitionally equal to f :⁴

$$f \equiv (\lambda x. f(x)).$$

This equality is the **uniqueness principle for function types**, because it shows that f is uniquely determined by its values.

The introduction of functions by definitions with explicit parameters can be reduced to simple definitions by using λ -abstraction: i.e., we can read a definition of $f : A \rightarrow B$ by

$$f(x) := \Phi$$

as

$$f := \lambda x. \Phi.$$

When doing calculations involving variables, we have to be careful when replacing a variable with an expression that also involves variables, because we want to preserve the binding structure of expressions. By the *binding structure* we mean the invisible link generated by binders such as λ , Π and Σ (the latter we are going to meet soon) between the place where the variable is introduced and where it is used. As an example, consider $f : \mathbb{N} \rightarrow (\mathbb{N} \rightarrow \mathbb{N})$ defined as

$$f(x) := \lambda y. x + y.$$

Now if we have assumed somewhere that $y : \mathbb{N}$, then what is $f(y)$? It would be wrong to just naively replace x by y everywhere in the expression “ $\lambda y. x + y$ ” defining $f(x)$, obtaining $\lambda y. y + y$, because this means that y gets **captured**. Previously, the substituted y was referring to our assumption, but now it is referring to the argument of the λ -abstraction. Hence, this naive substitution would destroy the binding structure, allowing us to perform calculations which are semantically unsound.

But what is $f(y)$ in this example? Note that bound (or “dummy”) variables such as y in the expression $\lambda y. x + y$ have only a local meaning, and can be consistently replaced by any other variable, preserving the binding structure. Indeed, $\lambda y. x + y$ is declared to be judgmentally equal⁵ to $\lambda z. x + z$. It follows that $f(y)$ is judgmentally equal to $\lambda z. y + z$, and that answers our question. (Instead of z , any variable distinct from y could have been used, yielding an equal result.)

Of course, this should all be familiar to any mathematician: it is the same phenomenon as the fact that if $f(x) := \int_1^2 \frac{dt}{x-t}$, then $f(t)$ is not $\int_1^2 \frac{dt}{t-t}$ but rather $\int_1^2 \frac{ds}{t-s}$. A λ -abstraction binds a dummy variable in exactly the same way that an integral does.

³Use of this equality is often referred to as **β -conversion** or **β -reduction**.

⁴Use of this equality is often referred to as **η -conversion** or **η -expansion**.

⁵Use of this equality is often referred to as **α -conversion**.

We have seen how to define functions in one variable. One way to define functions in several variables would be to use the cartesian product, which will be introduced later; a function with parameters A and B and results in C would be given the type $f : A \times B \rightarrow C$. However, there is another choice that avoids using product types, which is called **currying** (after the mathematician Haskell Curry).

The idea of currying is to represent a function of two inputs $a : A$ and $b : B$ as a function which takes *one* input $a : A$ and returns *another function*, which then takes a second input $b : B$ and returns the result. That is, we consider two-variable functions to belong to an iterated function type, $f : A \rightarrow (B \rightarrow C)$. We may also write this without the parentheses, as $f : A \rightarrow B \rightarrow C$, with associativity to the right as the default convention. Then given $a : A$ and $b : B$, we can apply f to a and then apply the result to b , obtaining $f(a)(b) : C$. To avoid the proliferation of parentheses, we allow ourselves to write $f(a)(b)$ as $f(a, b)$ even though there are no products involved. When omitting parentheses around function arguments entirely, we write $f\ a\ b$ for $(f\ a)\ b$, with the default associativity now being to the left so that f is applied to its arguments in the correct order.

Our notation for definitions with explicit parameters extends to this situation: we can define a named function $f : A \rightarrow B \rightarrow C$ by giving an equation

$$f(x, y) := \Phi$$

where $\Phi : C$ assuming $x : A$ and $y : B$. Using λ -abstraction this corresponds to

$$f := \lambda x. \lambda y. \Phi,$$

which may also be written as

$$f := x \mapsto y \mapsto \Phi.$$

We can also implicitly abstract over multiple variables by writing multiple blanks, e.g. $g(-, -)$ means $\lambda x. \lambda y. g(x, y)$. Currying a function of three or more arguments is a straightforward extension of what we have just described.

1.3 Universes and families

So far, we have been using the expression “ A is a type” informally. We are going to make this more precise by introducing **universes**. A universe is a type whose elements are types. As in naive set theory, we might wish for a universe of all types \mathcal{U}_∞ including itself (that is, with $\mathcal{U}_\infty : \mathcal{U}_\infty$). However, as in set theory, this is unsound, i.e. we can deduce from it that every type, including the empty type representing the proposition False (see §1.7), is inhabited. For instance, using a representation of sets as trees, we can directly encode Russell’s paradox [Coq92a].

To avoid the paradox we introduce a hierarchy of universes

$$\mathcal{U}_0 : \mathcal{U}_1 : \mathcal{U}_2 : \dots$$

where every universe \mathcal{U}_i is an element of the next universe \mathcal{U}_{i+1} . Moreover, we assume that our universes are **cumulative**, that is that all the elements of the i^{th} universe are also elements of the $(i + 1)^{\text{st}}$ universe, i.e. if $A : \mathcal{U}_i$ then also $A : \mathcal{U}_{i+1}$. This is convenient, but has the slightly unpleasant consequence that elements no longer have unique types, and is a bit tricky in other ways that need not concern us here; see the Notes.

When we say that A is a type, we mean that it inhabits some universe \mathcal{U}_i . We usually want to avoid mentioning the level i explicitly, and just assume that levels can be assigned in a consistent way; thus we may write $A : \mathcal{U}$ omitting the level. This way we can even write $\mathcal{U} : \mathcal{U}$, which can be read as $\mathcal{U}_i : \mathcal{U}_{i+1}$, having left the indices implicit. Writing universes in this style is referred to as **typical ambiguity**. It is convenient but a bit dangerous, since it allows us to write valid-looking proofs that reproduce the paradoxes of self-reference. If there is any doubt about whether an argument is correct, the way to check it is to try to assign levels consistently to all universes appearing in it. When some universe \mathcal{U} is assumed, we may refer to types belonging to \mathcal{U} as **small types**.

To model a collection of types varying over a given type A , we use functions $B : A \rightarrow \mathcal{U}$ whose codomain is a universe. These functions are called **families of types** (or sometimes *dependent types*); they correspond to families of sets as used in set theory.

An example of a type family is the family of finite sets $\text{Fin} : \mathbb{N} \rightarrow \mathcal{U}$, where $\text{Fin}(n)$ is a type with exactly n elements. (We cannot *define* the family Fin yet — indeed, we have not even introduced its domain \mathbb{N} yet — but we will be able to soon; see Exercise 1.9.) We may denote the elements of $\text{Fin}(n)$ by $0_n, 1_n, \dots, (n-1)_n$, with subscripts to emphasize that the elements of $\text{Fin}(n)$ are different from those of $\text{Fin}(m)$ if n is different from m , and all are different from the ordinary natural numbers (which we will introduce in §1.9).

A more trivial (but very important) example of a type family is the **constant** type family at a type $B : \mathcal{U}$, which is of course the constant function $(\lambda(x:A). B) : A \rightarrow \mathcal{U}$.

As a *non-example*, in our version of type theory there is no type family “ $\lambda(i:\mathbb{N}). \mathcal{U}_i$ ”. Indeed, there is no universe large enough to be its codomain. Moreover, we do not even identify the indices i of the universes \mathcal{U}_i with the natural numbers \mathbb{N} of type theory (the latter to be introduced in §1.9).

1.4 Dependent function types (Π -types)

In type theory we often use a more general version of function types, called a **Π -type** or **dependent function type**. The elements of a Π -type are functions whose codomain type can vary depending on the element of the domain to which the function is applied, called **dependent functions**. The name “ Π -type” is used because this type can also be regarded as the cartesian product over a given type.

Given a type $A : \mathcal{U}$ and a family $B : A \rightarrow \mathcal{U}$, we may construct the type of dependent functions $\prod_{(x:A)} B(x) : \mathcal{U}$. There are many alternative notations for this type, such as

$$\prod_{(x:A)} B(x) \quad \prod_{(x:A)} B(x) \quad \prod(x:A), B(x).$$

If B is a constant family, then the dependent product type is the ordinary function type:

$$\prod_{(x:A)} B \equiv (A \rightarrow B).$$

Indeed, all the constructions of Π -types are generalizations of the corresponding constructions on ordinary function types.

We can introduce dependent functions by explicit definitions: to define $f : \prod_{(x:A)} B(x)$, where f is the name of a dependent function to be defined, we need an expression $\Phi : B(x)$ possibly involving the variable $x : A$, and we write

$$f(x) := \Phi \quad \text{for } x : A.$$

Alternatively, we can use **λ -abstraction**

$$\lambda x. \Phi : \prod_{x:A} B(x). \quad (1.4.1)$$

As with non-dependent functions, we can **apply** a dependent function $f : \prod_{(x:A)} B(x)$ to an argument $a : A$ to obtain an element $f(a) : B(a)$. The equalities are the same as for the ordinary function type, i.e. we have the computation rule given $a : A$ we have $f(a) \equiv \Phi'$ and $(\lambda x. \Phi)(a) \equiv \Phi'$, where Φ' is obtained by replacing all occurrences of x in Φ by a (avoiding variable capture, as always). Similarly, we have the uniqueness principle $f \equiv (\lambda x. f(x))$ for any $f : \prod_{(x:A)} B(x)$.

As an example, recall from §1.3 that there is a type family $\text{Fin} : \mathbb{N} \rightarrow \mathcal{U}$ whose values are the standard finite sets, with elements $0_n, 1_n, \dots, (n-1)_n : \text{Fin}(n)$. There is then a dependent function $\text{fmax} : \prod_{(n:\mathbb{N})} \text{Fin}(n+1)$ which returns the “largest” element of each nonempty finite type, $\text{fmax}(n) := n_{n+1}$. As was the case for Fin itself, we cannot define fmax yet, but we will be able to soon; see Exercise 1.9.

Another important class of dependent function types, which we can define now, are functions which are **polymorphic** over a given universe. A polymorphic function is one which takes a type as one of its arguments, and then acts on elements of that type (or of other types constructed from it). An example is the polymorphic identity function $\text{id} : \prod_{(A:\mathcal{U})} A \rightarrow A$, which we define by $\text{id} := \lambda(A:\mathcal{U}). \lambda(x:A). x$. (Like λ -abstractions, Π s automatically scope over the rest of the expression unless delimited; thus $\text{id} : \prod_{(A:\mathcal{U})} A \rightarrow A$ means $\text{id} : \prod_{(A:\mathcal{U})}(A \rightarrow A)$. This convention, though unusual in mathematics, is common in type theory.)

We sometimes write some arguments of a dependent function as subscripts. For instance, we might equivalently define the polymorphic identity function by $\text{id}_A(x) := x$. Moreover, if an argument can be inferred from context, we may omit it altogether. For instance, if $a : A$, then writing $\text{id}(a)$ is unambiguous, since id must mean id_A in order for it to be applicable to a .

Another, less trivial, example of a polymorphic function is the “swap” operation that switches the order of the arguments of a (curried) two-argument function:

$$\text{swap} : \prod_{(A:\mathcal{U})} \prod_{(B:\mathcal{U})} \prod_{(C:\mathcal{U})} (A \rightarrow B \rightarrow C) \rightarrow (B \rightarrow A \rightarrow C).$$

We can define this by

$$\text{swap}(A, B, C, g) := \lambda b. \lambda a. g(a)(b).$$

We might also equivalently write the type arguments as subscripts:

$$\text{swap}_{A,B,C}(g)(b, a) := g(a, b).$$

Note that as we did for ordinary functions, we use currying to define dependent functions with several arguments (such as swap). However, in the dependent case the second domain may depend on the first one, and the codomain may depend on both. That is, given $A : \mathcal{U}$ and type families $B : A \rightarrow \mathcal{U}$ and $C : \prod_{(x:A)} B(x) \rightarrow \mathcal{U}$, we may construct the type $\prod_{(x:A)} \prod_{(y:B(x))} C(x, y)$ of functions with two arguments. In the case when B is constant and equal to A , we may condense the notation and write $\prod_{(x,y:A)}$; for instance, the type of swap could also be written as

$$\text{swap} : \prod_{A,B,C:\mathcal{U}} (A \rightarrow B \rightarrow C) \rightarrow (B \rightarrow A \rightarrow C).$$

Finally, given $f : \prod_{(x:A)} \prod_{(y:B(x))} C(x, y)$ and arguments $a : A$ and $b : B(a)$, we have $f(a)(b) : C(a, b)$, which, as before, we write as $f(a, b) : C(a, b)$.

1.5 Product types

Given types $A, B : \mathcal{U}$ we introduce the type $A \times B : \mathcal{U}$, which we call their **cartesian product**. We also introduce a nullary product type, called the **unit type** $\mathbf{1} : \mathcal{U}$. We intend the elements of $A \times B$ to be pairs $(a, b) : A \times B$, where $a : A$ and $b : B$, and the only element of $\mathbf{1}$ to be some particular object $\star : \mathbf{1}$. However, unlike in set theory, where we define ordered pairs to be particular sets and then collect them all together into the cartesian product, in type theory, ordered pairs are a primitive concept, as are functions.

Remark 1.5.1. There is a general pattern for introduction of a new kind of type in type theory. We have already seen this pattern in §§1.2 and 1.4⁶, so it is worth emphasizing the general form. To specify a type, we specify:

- (i) how to form new types of this kind, via **formation rules**. (For example, we can form the function type $A \rightarrow B$ when A is a type and when B is a type. We can form the dependent function type $\prod_{(x:A)} B(x)$ when A is a type and $B(x)$ is a type for $x : A$.)
- (ii) how to construct elements of that type. These are called the type's **constructors** or **introduction rules**. (For example, a function type has one constructor, λ -abstraction. Recall that a direct definition like $f(x) := 2x$ can equivalently be phrased as a λ -abstraction $f := \lambda x. 2x$.)
- (iii) how to use elements of that type. These are called the type's **eliminators** or **elimination rules**. (For example, the function type has one eliminator, namely function application.)
- (iv) a **computation rule**⁷, which expresses how an eliminator acts on a constructor. (For example, for functions, the computation rule states that $(\lambda x. \Phi)(a)$ is judgmentally equal to the substitution of a for x in Φ .)
- (v) an optional **uniqueness principle**⁸, which expresses uniqueness of maps into or out of that type. For some types, the uniqueness principle characterizes maps into the type, by stating that every element of the type is uniquely determined by the results of applying eliminators to it, and can be reconstructed from those results by applying a constructor—thus expressing how constructors act on eliminators, dually to the computation rule. (For example, for functions, the uniqueness principle says that any function f is judgmentally equal to the “expanded” function $\lambda x. f(x)$, and thus is uniquely determined by its values.) For other types, the uniqueness principle says that every map (function) *from* that type is uniquely determined by some data. (An example is the coproduct type introduced in §1.7, whose uniqueness principle is mentioned in §2.15.)

When the uniqueness principle is not taken as a rule of judgmental equality, it is often nevertheless provable as a *propositional* equality from the other rules for the type. In this case we call it a **propositional uniqueness principle**. (In later chapters we will also occasionally encounter *propositional computation rules*.)

The inference rules in ?? are organized and named accordingly; see, for example, ??, where each possibility is realized.

The way to construct pairs is obvious: given $a : A$ and $b : B$, we may form $(a, b) : A \times B$. Similarly, there is a unique way to construct elements of $\mathbf{1}$, namely we have $\star : \mathbf{1}$. We expect that

⁶The description of universes above is an exception.

⁷also referred to as β -reduction

⁸also referred to as η -expansion

“every element of $A \times B$ is a pair”, which is the uniqueness principle for products; we do not assert this as a rule of type theory, but we will prove it later on as a propositional equality.

Now, how can we *use* pairs, i.e. how can we define functions out of a product type? Let us first consider the definition of a non-dependent function $f : A \times B \rightarrow C$. Since we intend the only elements of $A \times B$ to be pairs, we expect to be able to define such a function by prescribing the result when f is applied to a pair (a, b) . We can prescribe these results by providing a function $g : A \rightarrow B \rightarrow C$. Thus, we introduce a new rule (the elimination rule for products), which says that for any such g , we can define a function $f : A \times B \rightarrow C$ by

$$f((a, b)) := g(a)(b).$$

We avoid writing $g(a, b)$ here, in order to emphasize that g is not a function on a product. (However, later on in the book we will often write $g(a, b)$ both for functions on a product and for curried functions of two variables.) This defining equation is the computation rule for product types.

Note that in set theory, we would justify the above definition of f by the fact that every element of $A \times B$ is an ordered pair, so that it suffices to define f on such pairs. By contrast, type theory reverses the situation: we assume that a function on $A \times B$ is well-defined as soon as we specify its values on pairs, and from this (or more precisely, from its more general version for dependent functions, below) we will be able to *prove* that every element of $A \times B$ is a pair. From a category-theoretic perspective, we can say that we define the product $A \times B$ to be left adjoint to the “exponential” $B \rightarrow C$, which we have already introduced.

As an example, we can derive the **projection** functions

$$\begin{aligned} \text{pr}_1 &: A \times B \rightarrow A \\ \text{pr}_2 &: A \times B \rightarrow B \end{aligned}$$

with the defining equations

$$\begin{aligned} \text{pr}_1((a, b)) &:= a \\ \text{pr}_2((a, b)) &:= b. \end{aligned}$$

Rather than invoking this principle of function definition every time we want to define a function, an alternative approach is to invoke it once, in a universal case, and then simply apply the resulting function in all other cases. That is, we may define a function of type

$$\text{rec}_{A \times B} : \prod_{C:U} (A \rightarrow B \rightarrow C) \rightarrow A \times B \rightarrow C \tag{1.5.2}$$

with the defining equation

$$\text{rec}_{A \times B}(C, g, (a, b)) := g(a)(b).$$

Then instead of defining functions such as pr_1 and pr_2 directly by a defining equation, we could define

$$\begin{aligned} \text{pr}_1 &:= \text{rec}_{A \times B}(A, \lambda a. \lambda b. a) \\ \text{pr}_2 &:= \text{rec}_{A \times B}(B, \lambda a. \lambda b. b). \end{aligned}$$

We refer to the function $\text{rec}_{A \times B}$ as the **recursor** for product types. The name “recursor” is a bit unfortunate here, since no recursion is taking place. It comes from the fact that product types

are a degenerate example of a general framework for inductive types, and for types such as the natural numbers, the recursor will actually be recursive. We may also speak of the **recursion principle** for cartesian products, meaning the fact that we can define a function $f : A \times B \rightarrow C$ as above by giving its value on pairs.

We leave it as a simple exercise to show that the recursor can be derived from the projections and vice versa.

We also have a recursor for the unit type:

$$\text{rec}_\mathbf{1} : \prod_{C:\mathcal{U}} C \rightarrow \mathbf{1} \rightarrow C$$

with the defining equation

$$\text{rec}_\mathbf{1}(C, c, \star) : \equiv c.$$

Although we include it to maintain the pattern of type definitions, the recursor for $\mathbf{1}$ is completely useless, because we could have defined such a function directly by simply ignoring the argument of type $\mathbf{1}$.

To be able to define *dependent* functions over the product type, we have to generalize the recursor. Given $C : A \times B \rightarrow \mathcal{U}$, we may define a function $f : \prod_{(x:A \times B)} C(x)$ by providing a function $g : \prod_{(x:A)} \prod_{(y:B)} C((x, y))$ with defining equation

$$f((x, y)) : \equiv g(x)(y).$$

For example, in this way we can prove the propositional uniqueness principle, which says that every element of $A \times B$ is equal to a pair. Specifically, we can construct a function

$$\text{uniq}_{A \times B} : \prod_{x:A \times B} ((\text{pr}_1(x), \text{pr}_2(x)) =_{A \times B} x).$$

Here we are using the identity type, which we are going to introduce below in §1.12. However, all we need to know now is that there is a reflexivity element $\text{refl}_x : x =_A x$ for any $x : A$. Given this, we can define

$$\text{uniq}_{A \times B}((a, b)) : \equiv \text{refl}_{(a,b)}.$$

This construction works, because in the case that $x : \equiv (a, b)$ we can calculate

$$(\text{pr}_1((a, b)), \text{pr}_2((a, b))) \equiv (a, b)$$

using the defining equations for the projections. Therefore,

$$\text{refl}_{(a,b)} : (\text{pr}_1((a, b)), \text{pr}_2((a, b))) = (a, b)$$

is well-typed, since both sides of the equality are judgmentally equal.

More generally, the ability to define dependent functions in this way means that to prove a property for all elements of a product, it is enough to prove it for its canonical elements, the ordered pairs. When we come to inductive types such as the natural numbers, the analogous property will be the ability to write proofs by induction. Thus, if we do as we did above and apply this principle once in the universal case, we call the resulting function **induction** for product types: given $A, B : \mathcal{U}$ we have

$$\text{ind}_{A \times B} : \prod_{C:A \times B \rightarrow \mathcal{U}} \left(\prod_{(x:A)} \prod_{(y:B)} C((x, y)) \right) \rightarrow \prod_{x:A \times B} C(x)$$

with the defining equation

$$\text{ind}_{A \times B}(C, g, (a, b)) : \equiv g(a)(b).$$

Similarly, we may speak of a dependent function defined on pairs being obtained from the **induction principle** of the cartesian product. It is easy to see that the recursor is just the special case of induction in the case that the family C is constant. Because induction describes how to use an element of the product type, induction is also called the **(dependent) eliminator**, and recursion the **non-dependent eliminator**.

Induction for the unit type turns out to be more useful than the recursor:

$$\text{ind}_1 : \prod_{C: 1 \rightarrow \mathcal{U}} C(\star) \rightarrow \prod_{x: 1} C(x)$$

with the defining equation

$$\text{ind}_1(C, c, \star) : \equiv c.$$

Induction enables us to prove the propositional uniqueness principle for 1 , which asserts that its only inhabitant is \star . That is, we can construct

$$\text{uniq}_1 : \prod_{x: 1} x = \star$$

by using the defining equations

$$\text{uniq}_1(\star) : \equiv \text{refl}_\star$$

or equivalently by using induction:

$$\text{uniq}_1 : \equiv \text{ind}_1(\lambda x. x = \star, \text{refl}_\star).$$

1.6 Dependent pair types (Σ -types)

Just as we generalized function types (§1.2) to dependent function types (§1.4), it is often useful to generalize the product types from §1.5 to allow the type of the second component of a pair to vary depending on the choice of the first component. This is called a **dependent pair type**, or **Σ -type**, because in set theory it corresponds to an indexed sum (in the sense of coproduct or disjoint union) over a given type.

Given a type $A : \mathcal{U}$ and a family $B : A \rightarrow \mathcal{U}$, the dependent pair type is written as $\sum_{(x:A)} B(x) : \mathcal{U}$. Alternative notations are

$$\sum_{(x:A)} B(x) \quad \sum_{(x:A)} B(x) \quad \sum_{(x:A)} B(x),$$

Like other binding constructs such as λ -abstractions and Π s, Σ s automatically scope over the rest of the expression unless delimited, so e.g. $\sum_{(x:A)} B(x) \times C(x)$ means $\sum_{(x:A)} (B(x) \times C(x))$.

The way to construct elements of a dependent pair type is by pairing: we have $(a, b) : \sum_{(x:A)} B(x)$ given $a : A$ and $b : B(a)$. If B is constant, then the dependent pair type is the ordinary cartesian product type:

$$\left(\sum_{x:A} B \right) \equiv (A \times B).$$

All the constructions on Σ -types arise as straightforward generalizations of the ones for product types, with dependent functions often replacing non-dependent ones.

For instance, the recursion principle says that to define a non-dependent function out of a Σ -type $f : (\sum_{(x:A)} B(x)) \rightarrow C$, we provide a function $g : \prod_{(x:A)} B(x) \rightarrow C$, and then we can define f via the defining equation

$$f((a, b)) := g(a)(b).$$

For instance, we can derive the first projection from a Σ -type:

$$\text{pr}_1 : \left(\sum_{x:A} B(x) \right) \rightarrow A$$

by the defining equation

$$\text{pr}_1((a, b)) := a.$$

However, since the type of the second component of a pair $(a, b) : \sum_{(x:A)} B(x)$ is $B(a)$, the second projection must be a *dependent* function, whose type involves the first projection function:

$$\text{pr}_2 : \prod_{p:\sum_{(x:A)} B(x)} B(\text{pr}_1(p)).$$

Thus we need the *induction* principle for Σ -types (the “dependent eliminator”). This says that to construct a dependent function out of a Σ -type into a family $C : (\sum_{(x:A)} B(x)) \rightarrow \mathcal{U}$, we need a function

$$g : \prod_{(a:A)} \prod_{(b:B(a))} C((a, b)).$$

We can then derive a function

$$f : \prod_{p:\sum_{(x:A)} B(x)} C(p)$$

with defining equation

$$f((a, b)) := g(a)(b).$$

Applying this with $C(p) := B(\text{pr}_1(p))$, we can define $\text{pr}_2 : \prod_{(p:\sum_{(x:A)} B(x))} B(\text{pr}_1(p))$ with the obvious equation

$$\text{pr}_2((a, b)) := b.$$

To convince ourselves that this is correct, we note that $B(\text{pr}_1((a, b))) \equiv B(a)$, using the defining equation for pr_1 , and indeed $b : B(a)$.

We can package the recursion and induction principles into the recursor for Σ :

$$\text{rec}_{\sum_{(x:A)} B(x)} : \prod_{(C:\mathcal{U})} \left(\prod_{(x:A)} B(x) \rightarrow C \right) \rightarrow \left(\sum_{(x:A)} B(x) \right) \rightarrow C$$

with the defining equation

$$\text{rec}_{\sum_{(x:A)} B(x)}(C, g, (a, b)) := g(a)(b)$$

and the corresponding induction operator:

$$\text{ind}_{\sum_{(x:A)} B(x)} : \prod_{(C:(\sum_{(x:A)} B(x)) \rightarrow \mathcal{U})} \left(\prod_{(a:A)} \prod_{(b:B(a))} C((a, b)) \right) \rightarrow \prod_{(p:\sum_{(x:A)} B(x))} C(p)$$

with the defining equation

$$\text{ind}_{\sum_{(x:A)} B(x)}(C, g, (a, b)) \equiv g(a)(b).$$

As before, the recursor is the special case of induction when the family C is constant.

As a further example, consider the following principle, where A and B are types and $R : A \rightarrow B \rightarrow \mathcal{U}$:

$$\text{ac} : \left(\prod_{(x:A)} \sum_{(y:B)} R(x, y) \right) \rightarrow \left(\sum_{(f:A \rightarrow B)} \prod_{(x:A)} R(x, f(x)) \right).$$

We may regard R as a “proof-relevant relation” between A and B , with $R(a, b)$ the type of witnesses for relatedness of $a : A$ and $b : B$. Then ac says intuitively that if we have a dependent function g assigning to every $a : A$ a dependent pair (b, r) where $b : B$ and $r : R(a, b)$, then we have a function $f : A \rightarrow B$ and a dependent function assigning to every $a : A$ a witness that $R(a, f(a))$. Our intuition tells us that we can just split up the values of g into their components. Indeed, using the projections we have just defined, we can define:

$$\text{ac}(g) \equiv \left(\lambda x. \text{pr}_1(g(x)), \lambda x. \text{pr}_2(g(x)) \right).$$

To verify that this is well-typed, note that if $g : \prod_{(x:A)} \sum_{(y:B)} R(x, y)$, we have

$$\begin{aligned} \lambda x. \text{pr}_1(g(x)) &: A \rightarrow B, \\ \lambda x. \text{pr}_2(g(x)) &: \prod_{(x:A)} R(x, \text{pr}_1(g(x))). \end{aligned}$$

Moreover, the type $\prod_{(x:A)} R(x, \text{pr}_1(g(x)))$ is the result of applying the type family $\lambda f. \prod_{(x:A)} R(x, f(x))$ being summed over in the codomain of ac to the function $\lambda x. \text{pr}_1(g(x))$:

$$\prod_{(x:A)} R(x, \text{pr}_1(g(x))) \equiv \left(\lambda f. \prod_{(x:A)} R(x, f(x)) \right) (\lambda x. \text{pr}_1(g(x))).$$

Thus, we have

$$\left(\lambda x. \text{pr}_1(g(x)), \lambda x. \text{pr}_2(g(x)) \right) : \sum_{(f:A \rightarrow B)} \prod_{(x:A)} R(x, f(x))$$

as required.

If we read Π as “for all” and Σ as “there exists”, then the type of the function ac expresses: *if for all $x : A$ there is a $y : B$ such that $R(x, y)$, then there is a function $f : A \rightarrow B$ such that for all $x : A$ we have $R(x, f(x))$* . Since this sounds like a version of the axiom of choice, the function ac has traditionally been called the **type-theoretic axiom of choice**, and as we have just shown, it can be proved directly from the rules of type theory, rather than having to be taken as an axiom. However, note that no choice is actually involved, since the choices have already been given to us in the premise: all we have to do is take it apart into two functions: one representing the choice and the other its correctness. In §3.8 we will give another formulation of an “axiom of choice” which is closer to the usual one.

Dependent pair types are often used to define types of mathematical structures, which commonly consist of several dependent pieces of data. To take a simple example, suppose we want to define a **magma** to be a type A together with a binary operation $m : A \rightarrow A \rightarrow A$. The precise meaning of the phrase “together with” (and the synonymous “equipped with”) is that “a magma” is a *pair* (A, m) consisting of a type $A : \mathcal{U}$ and an operation $m : A \rightarrow A \rightarrow A$. Since the type $A \rightarrow A \rightarrow A$ of the second component m of this pair depends on its first component A ,

such pairs belong to a dependent pair type. Thus, the definition “a magma is a type A together with a binary operation $m : A \rightarrow A \rightarrow A$ ” should be read as defining *the type of magmas* to be

$$\text{Magma} := \sum_{A:\mathcal{U}} (A \rightarrow A \rightarrow A).$$

Given a magma, we extract its underlying type (its “carrier”) with the first projection pr_1 , and its operation with the second projection pr_2 . Of course, structures built from more than two pieces of data require iterated pair types, which may be only partially dependent; for instance the type of pointed magmas ($\text{magmas } (A, m)$ equipped with a basepoint $e : A$) is

$$\text{PointedMagma} := \sum_{A:\mathcal{U}} (A \rightarrow A \rightarrow A) \times A.$$

We generally also want to impose axioms on such a structure, e.g. to make a pointed magma into a monoid or a group. This can also be done using Σ -types; see §1.11.

In the rest of the book, we will sometimes make definitions of this sort explicit, but eventually we trust the reader to translate them from English into Σ -types. We also generally follow the common mathematical practice of using the same letter for a structure of this sort and for its carrier (which amounts to leaving the appropriate projection function implicit in the notation): that is, we will speak of a magma A with its operation $m : A \rightarrow A \rightarrow A$.

Note that the canonical elements of PointedMagma are of the form $(A, (m, e))$ where $A : \mathcal{U}$, $m : A \rightarrow A \rightarrow A$, and $e : A$. Because of the frequency with which iterated Σ -types of this sort arise, we use the usual notation of ordered triples, quadruples and so on to stand for nested pairs (possibly dependent) associating to the right. That is, we have $(x, y, z) := (x, (y, z))$ and $(x, y, z, w) := (x, (y, (z, w)))$, etc.

1.7 Coproduct types

Given $A, B : \mathcal{U}$, we introduce their **coproduct** type $A + B : \mathcal{U}$. This corresponds to the *disjoint union* in set theory, and we may also use that name for it. In type theory, as was the case with functions and products, the coproduct must be a fundamental construction, since there is no previously given notion of “union of types”. We also introduce a nullary version: the **empty type** $\mathbf{0} : \mathcal{U}$.

There are two ways to construct elements of $A + B$, either as $\text{inl}(a) : A + B$ for $a : A$, or as $\text{inr}(b) : A + B$ for $b : B$. (The names inl and inr are short for “left injection” and “right injection”.) There are no ways to construct elements of the empty type.

To construct a non-dependent function $f : A + B \rightarrow C$, we need functions $g_0 : A \rightarrow C$ and $g_1 : B \rightarrow C$. Then f is defined via the defining equations

$$\begin{aligned} f(\text{inl}(a)) &:= g_0(a), \\ f(\text{inr}(b)) &:= g_1(b). \end{aligned}$$

That is, the function f is defined by **case analysis**. As before, we can derive the recursor:

$$\text{rec}_{A+B} : \prod_{(C:\mathcal{U})} (A \rightarrow C) \rightarrow (B \rightarrow C) \rightarrow A + B \rightarrow C$$

with the defining equations

$$\begin{aligned} \text{rec}_{A+B}(C, g_0, g_1, \text{inl}(a)) &:= g_0(a), \\ \text{rec}_{A+B}(C, g_0, g_1, \text{inr}(b)) &:= g_1(b). \end{aligned}$$

We can always construct a function $f : \mathbf{0} \rightarrow C$ without having to give any defining equations, because there are no elements of $\mathbf{0}$ on which to define f . Thus, the recursor for $\mathbf{0}$ is

$$\text{rec}_0 : \prod_{(C:\mathcal{U})} \mathbf{0} \rightarrow C,$$

which constructs the canonical function from the empty type to any other type. Logically, it corresponds to the principle *ex falso quodlibet*.

To construct a dependent function $f : \prod_{(x:A+B)} C(x)$ out of a coproduct, we assume as given the family $C : (A + B) \rightarrow \mathcal{U}$, and require

$$\begin{aligned} g_0 &: \prod_{a:A} C(\text{inl}(a)), \\ g_1 &: \prod_{b:B} C(\text{inr}(b)). \end{aligned}$$

This yields f with the defining equations:

$$\begin{aligned} f(\text{inl}(a)) &:= g_0(a), \\ f(\text{inr}(b)) &:= g_1(b). \end{aligned}$$

We package this scheme into the induction principle for coproducts:

$$\text{ind}_{A+B} : \prod_{(C:(A+B) \rightarrow \mathcal{U})} \left(\prod_{(a:A)} C(\text{inl}(a)) \right) \rightarrow \left(\prod_{(b:B)} C(\text{inr}(b)) \right) \rightarrow \prod_{(x:A+B)} C(x).$$

As before, the recursor arises in the case that the family C is constant.

The induction principle for the empty type

$$\text{ind}_0 : \prod_{(C:\mathbf{0} \rightarrow \mathcal{U})} \prod_{(z:\mathbf{0})} C(z)$$

gives us a way to define a trivial dependent function out of the empty type.

1.8 The type of booleans

The type of booleans $\mathbf{2} : \mathcal{U}$ is intended to have exactly two elements $0_2, 1_2 : \mathbf{2}$. It is clear that we could construct this type out of coproduct and unit types as $\mathbf{1} + \mathbf{1}$. However, since it is used frequently, we give the explicit rules here. Indeed, we are going to observe that we can also go the other way and derive binary coproducts from Σ -types and $\mathbf{2}$.

To derive a function $f : \mathbf{2} \rightarrow C$ we need $c_0, c_1 : C$ and add the defining equations

$$\begin{aligned} f(0_2) &:= c_0, \\ f(1_2) &:= c_1. \end{aligned}$$

The recursor corresponds to the if-then-else construct in functional programming:

$$\text{rec}_2 : \prod_{C:\mathcal{U}} C \rightarrow C \rightarrow \mathbf{2} \rightarrow C$$

with the defining equations

$$\begin{aligned} \text{rec}_2(C, c_0, c_1, 0_2) &:= c_0, \\ \text{rec}_2(C, c_0, c_1, 1_2) &:= c_1. \end{aligned}$$

Given $C : \mathbf{2} \rightarrow \mathcal{U}$, to derive a dependent function $f : \prod_{(x:\mathbf{2})} C(x)$ we need $c_0 : C(0_2)$ and $c_1 : C(1_2)$, in which case we can give the defining equations

$$\begin{aligned} f(0_2) &:= c_0, \\ f(1_2) &:= c_1. \end{aligned}$$

We package this up into the induction principle

$$\text{ind}_2 : \prod_{(C:\mathbf{2} \rightarrow \mathcal{U})} C(0_2) \rightarrow C(1_2) \rightarrow \prod_{(x:\mathbf{2})} C(x)$$

with the defining equations

$$\begin{aligned} \text{ind}_2(C, c_0, c_1, 0_2) &:= c_0, \\ \text{ind}_2(C, c_0, c_1, 1_2) &:= c_1. \end{aligned}$$

As an example, using the induction principle we can deduce that, as we expect, every element of $\mathbf{2}$ is either 1_2 or 0_2 . As before, in order to state this we use the equality types which we have not yet introduced, but we need only the fact that everything is equal to itself by $\text{refl}_x : x = x$. Thus, we construct an element of

$$\prod_{x:\mathbf{2}} (x = 0_2) + (x = 1_2), \tag{1.8.1}$$

i.e. a function assigning to each $x : \mathbf{2}$ either an equality $x = 0_2$ or an equality $x = 1_2$. We define this element using the induction principle for $\mathbf{2}$, with $C(x) := (x = 0_2) + (x = 1_2)$; the two inputs are $\text{inl}(\text{refl}_{0_2}) : C(0_2)$ and $\text{inr}(\text{refl}_{1_2}) : C(1_2)$. In other words, our element of (1.8.1) is

$$\text{ind}_2(\lambda x. (x = 0_2) + (x = 1_2), \text{inl}(\text{refl}_{0_2}), \text{inr}(\text{refl}_{1_2})).$$

We have remarked that Σ -types can be regarded as analogous to indexed disjoint unions, while coproducts are binary disjoint unions. It is natural to expect that a binary disjoint union $A + B$ could be constructed as an indexed one over the two-element type $\mathbf{2}$. For this we need a type family $P : \mathbf{2} \rightarrow \mathcal{U}$ such that $P(0_2) \equiv A$ and $P(1_2) \equiv B$. Indeed, we can obtain such a family precisely by the recursion principle for $\mathbf{2}$. (The ability to define *type families* by induction and recursion, using the fact that the universe \mathcal{U} is itself a type, is a subtle and important aspect of type theory.) Thus, we could have defined

$$A + B := \sum_{x:\mathbf{2}} \text{rec}_2(\mathcal{U}, A, B, x)$$

with

$$\begin{aligned} \text{inl}(a) &:= (0_2, a), \\ \text{inr}(b) &:= (1_2, b). \end{aligned}$$

We leave it as an exercise to derive the induction principle of a coproduct type from this definition. (See also Exercise 1.5 and §5.2.)

We can apply the same idea to products and Π -types: we could have defined

$$A \times B := \prod_{x:\mathbf{2}} \text{rec}_2(\mathcal{U}, A, B, x).$$

Pairs could then be constructed using induction for $\mathbf{2}$:

$$(a, b) := \text{ind}_2(\text{rec}_2(\mathcal{U}, A, B), a, b)$$

while the projections are straightforward applications

$$\begin{aligned}\text{pr}_1(p) &:= p(0_2), \\ \text{pr}_2(p) &:= p(1_2).\end{aligned}$$

The derivation of the induction principle for binary products defined in this way is a bit more involved, and requires function extensionality, which we will introduce in §2.9. Moreover, we do not get the same judgmental equalities; see Exercise 1.6. This is a recurrent issue when encoding one type as another; we will return to it in §5.5.

We may occasionally refer to the elements 0_2 and 1_2 of $\mathbf{2}$ as “false” and “true” respectively. However, note that unlike in classical mathematics, we do not use elements of $\mathbf{2}$ as truth values or as propositions. (Instead we identify propositions with types; see §1.11.) In particular, the type $A \rightarrow \mathbf{2}$ is not generally the power set of A ; it represents only the “decidable” subsets of A (see Chapter 3).

1.9 The natural numbers

So far we have rules for constructing new types by abstract operations, but for doing concrete mathematics we also require some concrete types, such as types of numbers. The most basic such is the type $\mathbb{N} : \mathcal{U}$ of natural numbers; once we have this we can construct integers, rational numbers, real numbers, and so on (see ??).

The elements of \mathbb{N} are constructed using $0 : \mathbb{N}$ and the successor operation $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$. When denoting natural numbers, we adopt the usual decimal notation $1 := \text{succ}(0)$, $2 := \text{succ}(1)$, $3 := \text{succ}(2)$,

The essential property of the natural numbers is that we can define functions by recursion and perform proofs by induction — where now the words “recursion” and “induction” have a more familiar meaning. To construct a non-dependent function $f : \mathbb{N} \rightarrow C$ out of the natural numbers by recursion, it is enough to provide a starting point $c_0 : C$ and a “next step” function $c_s : \mathbb{N} \rightarrow C \rightarrow C$. This gives rise to f with the defining equations

$$\begin{aligned}f(0) &:= c_0, \\ f(\text{succ}(n)) &:= c_s(n, f(n)).\end{aligned}$$

We say that f is defined by **primitive recursion**.

As an example, we look at how to define a function on natural numbers which doubles its argument. In this case we have $C := \mathbb{N}$. We first need to supply the value of $\text{double}(0)$, which is easy: we put $c_0 := 0$. Next, to compute the value of $\text{double}(\text{succ}(n))$ for a natural number n , we first compute the value of $\text{double}(n)$ and then perform the successor operation twice. This is captured by the recurrence $c_s(n, y) := \text{succ}(\text{succ}(y))$. Note that the second argument y of c_s stands for the result of the *recursive call* $\text{double}(n)$.

Defining $\text{double} : \mathbb{N} \rightarrow \mathbb{N}$ by primitive recursion in this way, therefore, we obtain the defining equations:

$$\begin{aligned}\text{double}(0) &:= 0 \\ \text{double}(\text{succ}(n)) &:= \text{succ}(\text{succ}(\text{double}(n))).\end{aligned}$$

This indeed has the correct computational behavior: for example, we have

$$\begin{aligned}
 \text{double}(2) &\equiv \text{double}(\text{succ}(\text{succ}(0))) \\
 &\equiv c_s(\text{succ}(0), \text{double}(\text{succ}(0))) \\
 &\equiv \text{succ}(\text{succ}(\text{double}(\text{succ}(0)))) \\
 &\equiv \text{succ}(\text{succ}(c_s(0, \text{double}(0)))) \\
 &\equiv \text{succ}(\text{succ}(\text{succ}(\text{succ}(\text{double}(0))))) \\
 &\equiv \text{succ}(\text{succ}(\text{succ}(\text{succ}(c_0)))) \\
 &\equiv \text{succ}(\text{succ}(\text{succ}(\text{succ}(0)))) \\
 &\equiv 4.
 \end{aligned}$$

We can define multi-variable functions by primitive recursion as well, by currying and allowing C to be a function type. For example, we define addition $\text{add} : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$ with $C := \mathbb{N} \rightarrow \mathbb{N}$ and the following “starting point” and “next step” data:

$$\begin{aligned}
 c_0 &: \mathbb{N} \rightarrow \mathbb{N} \\
 c_0(n) &:= n \\
 c_s &: \mathbb{N} \rightarrow (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow (\mathbb{N} \rightarrow \mathbb{N}) \\
 c_s(m, g)(n) &:= \text{succ}(g(n)).
 \end{aligned}$$

We thus obtain $\text{add} : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$ satisfying the definitional equalities

$$\begin{aligned}
 \text{add}(0, n) &\equiv n \\
 \text{add}(\text{succ}(m), n) &\equiv \text{succ}(\text{add}(m, n)).
 \end{aligned}$$

As usual, we write $\text{add}(m, n)$ as $m + n$. The reader is invited to verify that $2 + 2 \equiv 4$.

As in previous cases, we can package the principle of primitive recursion into a recursor:

$$\text{rec}_{\mathbb{N}} : \prod_{(C:\mathcal{U})} C \rightarrow (\mathbb{N} \rightarrow C \rightarrow C) \rightarrow \mathbb{N} \rightarrow C$$

with the defining equations

$$\begin{aligned}
 \text{rec}_{\mathbb{N}}(C, c_0, c_s, 0) &:= c_0, \\
 \text{rec}_{\mathbb{N}}(C, c_0, c_s, \text{succ}(n)) &:= c_s(n, \text{rec}_{\mathbb{N}}(C, c_0, c_s, n)).
 \end{aligned}$$

Using $\text{rec}_{\mathbb{N}}$ we can present double and add as follows:

$$\text{double} := \text{rec}_{\mathbb{N}}(\mathbb{N}, 0, \lambda n. \lambda y. \text{succ}(\text{succ}(y))) \tag{1.9.1}$$

$$\text{add} := \text{rec}_{\mathbb{N}}(\mathbb{N} \rightarrow \mathbb{N}, \lambda n. n, \lambda n. \lambda g. \lambda m. \text{succ}(g(m))). \tag{1.9.2}$$

Of course, all functions definable only using the primitive recursion principle will be *computable*. (The presence of higher function types — that is, functions with other functions as arguments — does, however, mean we can define more than the usual primitive recursive functions; see e.g. Exercise 1.10.) This is appropriate in constructive mathematics; in §§3.4 and 3.8 we will see how to augment type theory so that we can define more general mathematical functions.

We now follow the same approach as for other types, generalizing primitive recursion to dependent functions to obtain an *induction principle*. Thus, assume as given a family $C : \mathbb{N} \rightarrow \mathcal{U}$,

an element $c_0 : C(0)$, and a function $c_s : \prod_{(n:\mathbb{N})} C(n) \rightarrow C(\text{succ}(n))$; then we can construct $f : \prod_{(n:\mathbb{N})} C(n)$ with the defining equations:

$$\begin{aligned} f(0) &:= c_0, \\ f(\text{succ}(n)) &:= c_s(n, f(n)). \end{aligned}$$

We can also package this into a single function

$$\text{ind}_{\mathbb{N}} : \prod_{(C:\mathbb{N} \rightarrow \mathcal{U})} C(0) \rightarrow \left(\prod_{(n:\mathbb{N})} C(n) \rightarrow C(\text{succ}(n)) \right) \rightarrow \prod_{(n:\mathbb{N})} C(n)$$

with the defining equations

$$\begin{aligned} \text{ind}_{\mathbb{N}}(C, c_0, c_s, 0) &:= c_0, \\ \text{ind}_{\mathbb{N}}(C, c_0, c_s, \text{succ}(n)) &:= c_s(n, \text{ind}_{\mathbb{N}}(C, c_0, c_s, n)). \end{aligned}$$

Here we finally see the connection to the classical notion of proof by induction. Recall that in type theory we represent propositions by types, and proving a proposition by inhabiting the corresponding type. In particular, a *property* of natural numbers is represented by a family of types $P : \mathbb{N} \rightarrow \mathcal{U}$. From this point of view, the above induction principle says that if we can prove $P(0)$, and if for any n we can prove $P(\text{succ}(n))$ assuming $P(n)$, then we have $P(n)$ for all n . This is, of course, exactly the usual principle of proof by induction on natural numbers.

As an example, consider how we might represent an explicit proof that $+$ is associative. (We will not actually write out proofs in this style, but it serves as a useful example for understanding how induction is represented formally in type theory.) To derive

$$\text{assoc} : \prod_{i,j,k:\mathbb{N}} i + (j + k) = (i + j) + k,$$

it is sufficient to supply

$$\text{assoc}_0 : \prod_{j,k:\mathbb{N}} 0 + (j + k) = (0 + j) + k$$

and

$$\text{assoc}_s : \prod_{i:\mathbb{N}} \left(\prod_{j,k:\mathbb{N}} i + (j + k) = (i + j) + k \right) \rightarrow \prod_{j,k:\mathbb{N}} \text{succ}(i) + (j + k) = (\text{succ}(i) + j) + k.$$

To derive assoc_0 , recall that $0 + n \equiv n$, and hence $0 + (j + k) \equiv j + k \equiv (0 + j) + k$. Hence we can just set

$$\text{assoc}_0(j, k) := \text{refl}_{j+k}.$$

For assoc_s , recall that the definition of $+$ gives $\text{succ}(m) + n \equiv \text{succ}(m + n)$, and hence

$$\begin{aligned} \text{succ}(i) + (j + k) &\equiv \text{succ}(i + (j + k)) \quad \text{and} \\ (\text{succ}(i) + j) + k &\equiv \text{succ}((i + j) + k). \end{aligned}$$

Thus, the output type of assoc_s is equivalently $\text{succ}(i + (j + k)) = \text{succ}((i + j) + k)$. But its input (the “inductive hypothesis”) yields $i + (j + k) = (i + j) + k$, so it suffices to invoke the fact that if two natural numbers are equal, then so are their successors. (We will prove this obvious fact in Lemma 2.2.1, using the induction principle of identity types.) We call this latter fact $\text{ap}_{\text{succ}} : (m =_{\mathbb{N}} n) \rightarrow (\text{succ}(m) =_{\mathbb{N}} \text{succ}(n))$, so we can define

$$\text{assoc}_s(i, h, j, k) := \text{ap}_{\text{succ}}(h(j, k)).$$

Putting these together with $\text{ind}_{\mathbb{N}}$, we obtain a proof of associativity.

1.10 Pattern matching and recursion

The natural numbers introduce an additional subtlety over the types considered up until now. In the case of coproducts, for instance, we could define a function $f : A + B \rightarrow C$ either with the recursor:

$$f := \text{rec}_{A+B}(C, g_0, g_1)$$

or by giving the defining equations:

$$\begin{aligned} f(\text{inl}(a)) &:= g_0(a) \\ f(\text{inr}(b)) &:= g_1(b). \end{aligned}$$

To go from the former expression of f to the latter, we simply use the computation rules for the recursor. Conversely, given any defining equations

$$\begin{aligned} f(\text{inl}(a)) &:= \Phi_0 \\ f(\text{inr}(b)) &:= \Phi_1 \end{aligned}$$

where Φ_0 and Φ_1 are expressions that may involve the variables a and b respectively, we can express these equations equivalently in terms of the recursor by using λ -abstraction:

$$f := \text{rec}_{A+B}(C, \lambda a. \Phi_0, \lambda b. \Phi_1).$$

In the case of the natural numbers, however, the “defining equations” of a function such as `double`:

$$\text{double}(0) := 0 \tag{1.10.1}$$

$$\text{double}(\text{succ}(n)) := \text{succ}(\text{succ}(\text{double}(n))) \tag{1.10.2}$$

involve *the function double itself* on the right-hand side. However, we would still like to be able to give these equations, rather than (1.9.1), as the definition of `double`, since they are much more convenient and readable. The solution is to read the expression “`double(n)`” on the right-hand side of (1.10.2) as standing in for the result of the recursive call, which in a definition of the form $\text{double} := \text{rec}_{\mathbb{N}}(\mathbb{N}, c_0, c_s)$ would be the second argument of c_s .

More generally, if we have a “definition” of a function $f : \mathbb{N} \rightarrow C$ such as

$$\begin{aligned} f(0) &:= \Phi_0 \\ f(\text{succ}(n)) &:= \Phi_s \end{aligned}$$

where Φ_0 is an expression of type C , and Φ_s is an expression of type C which may involve the variable n and also the symbol “ $f(n)$ ”, we may translate it to a definition

$$f := \text{rec}_{\mathbb{N}}(C, \Phi_0, \lambda n. \lambda r. \Phi'_s)$$

where Φ'_s is obtained from Φ_s by replacing all occurrences of “ $f(n)$ ” by the new variable r .

This style of defining functions by recursion (or, more generally, dependent functions by induction) is so convenient that we frequently adopt it. It is called **definition by pattern matching**. Of course, it is very similar to how a computer programmer may define a recursive function with a body that literally contains recursive calls to itself. However, unlike the programmer, we are restricted in what sort of recursive calls we can make: in order for such a definition to be

re-expressible using the recursion principle, the function f being defined can only appear in the body of $f(\text{succ}(n))$ as part of the composite symbol “ $f(n)$ ”. Otherwise, we could write nonsense functions such as

$$\begin{aligned} f(0) &:= 0 \\ f(\text{succ}(n)) &:= f(\text{succ}(\text{succ}(n))). \end{aligned}$$

If a programmer wrote such a function, it would simply call itself forever on any positive input, going into an infinite loop and never returning a value. In mathematics, however, to be worthy of the name, a *function* must always associate a unique output value to every input value, so this would be unacceptable.

This point will be even more important when we introduce more complicated inductive types in Chapters 5 and 6 and ???. Whenever we introduce a new kind of inductive definition, we always begin by deriving its induction principle. Only then do we introduce an appropriate sort of “pattern matching” which can be justified as a shorthand for the induction principle.

1.11 Propositions as types

As mentioned in the introduction, to show that a proposition is true in type theory corresponds to exhibiting an element of the type corresponding to that proposition. We regard the elements of this type as *evidence* or *witnesses* that the proposition is true. (They are sometimes even called *proofs*, but this terminology can be misleading, so we generally avoid it.) In general, however, we will not construct witnesses explicitly; instead we present the proofs in ordinary mathematical prose, in such a way that they could be translated into an element of a type. This is no different from reasoning in classical set theory, where we don’t expect to see an explicit derivation using the rules of predicate logic and the axioms of set theory.

However, the type-theoretic perspective on proofs is nevertheless different in important ways. The basic principle of the logic of type theory is that a proposition is not merely true or false, but rather can be seen as the collection of all possible witnesses of its truth. Under this conception, proofs are not just the means by which mathematics is communicated, but rather are mathematical objects in their own right, on a par with more familiar objects such as numbers, mappings, groups, and so on. Thus, since types classify the available mathematical objects and govern how they interact, propositions are nothing but special types — namely, types whose elements are proofs.

The basic observation which makes this identification feasible is that we have the following natural correspondence between *logical* operations on propositions, expressed in English, and *type-theoretic* operations on their corresponding types of witnesses.

English	Type Theory
True	1
False	0
A and B	$A \times B$
A or B	$A + B$
If A then B	$A \rightarrow B$
A if and only if B	$(A \rightarrow B) \times (B \rightarrow A)$
Not A	$A \rightarrow \mathbf{0}$

The point of the correspondence is that in each case, the rules for constructing and using elements of the type on the right correspond to the rules for reasoning about the proposition on the left. For instance, the basic way to prove a statement of the form “ A and B ” is to prove A and also prove B , while the basic way to construct an element of $A \times B$ is as a pair (a, b) , where a is an element (or witness) of A and b is an element (or witness) of B . And if we want to use “ A and B ” to prove something else, we are free to use both A and B in doing so, analogously to how the induction principle for $A \times B$ allows us to construct a function out of it by using elements of A and of B .

Similarly, the basic way to prove an implication “if A then B ” is to assume A and prove B , while the basic way to construct an element of $A \rightarrow B$ is to give an expression which denotes an element (witness) of B which may involve an unspecified variable element (witness) of type A . And the basic way to use an implication “if A then B ” is deduce B if we know A , analogously to how we can apply a function $f : A \rightarrow B$ to an element of A to produce an element of B . We strongly encourage the reader to do the exercise of verifying that the rules governing the other type constructors translate sensibly into logic.

Of special note is that the empty type **0** corresponds to falsity. When speaking logically, we refer to an inhabitant of **0** as a **contradiction**: thus there is no way to prove a contradiction,⁹ while from a contradiction anything can be derived. We also define the **negation** of a type A as

$$\neg A := A \rightarrow \mathbf{0}.$$

Thus, a witness of $\neg A$ is a function $A \rightarrow \mathbf{0}$, which we may construct by assuming $x : A$ and deriving an element of **0**. Note that although the logic we obtain is “constructive”, as discussed in the introduction, this sort of “proof by contradiction” (assume A and derive a contradiction, concluding $\neg A$) is perfectly valid constructively: it is simply invoking the *meaning* of “negation”. The sort of “proof by contradiction” which is disallowed is to assume $\neg A$ and derive a contradiction as a way of proving A . Constructively, such an argument would only allow us to conclude $\neg\neg A$, and the reader can verify that there is no obvious way to get from $\neg\neg A$ (that is, from $(A \rightarrow \mathbf{0}) \rightarrow \mathbf{0}$) to A .

The above translation of logical connectives into type-forming operations is referred to as **propositions as types**: it gives us a way to translate propositions and their proofs, written in English, into types and their elements. For example, suppose we want to prove the following tautology (one of “de Morgan’s laws”):

$$\text{"If not } A \text{ and not } B, \text{ then not } (A \text{ or } B)". \quad (1.11.1)$$

An ordinary English proof of this fact might go as follows.

Suppose not A and not B , and also suppose A or B ; we will derive a contradiction. There are two cases. If A holds, then since not A , we have a contradiction. Similarly, if B holds, then since not B , we also have a contradiction. Thus we have a contradiction in either case, so not (A or B).

Now, the type corresponding to our tautology (1.11.1), according to the rules given above, is

$$(A \rightarrow \mathbf{0}) \times (B \rightarrow \mathbf{0}) \rightarrow (A + B \rightarrow \mathbf{0}) \quad (1.11.2)$$

⁹More precisely, there is no *basic* way to prove a contradiction, i.e. **0** has no constructors. If our type theory were inconsistent, then there would be some more complicated way to construct an element of **0**.

so we should be able to translate the above proof into an element of this type.

As an example of how such a translation works, let us describe how a mathematician reading the above English proof might simultaneously construct, in his or her head, an element of (1.11.2). The introductory phrase “Suppose not A and not B ” translates into defining a function, with an implicit application of the recursion principle for the cartesian product in its domain $(A \rightarrow \mathbf{0}) \times (B \rightarrow \mathbf{0})$. This introduces unnamed variables (hypotheses) of types $A \rightarrow \mathbf{0}$ and $B \rightarrow \mathbf{0}$. When translating into type theory, we have to give these variables names; let us call them x and y . At this point our partial definition of an element of (1.11.2) can be written as

$$f((x, y)) := \square : A + B \rightarrow \mathbf{0}$$

with a “hole” \square of type $A + B \rightarrow \mathbf{0}$ indicating what remains to be done. (We could equivalently write $f := \text{rec}_{(A \rightarrow \mathbf{0}) \times (B \rightarrow \mathbf{0})}(A + B \rightarrow \mathbf{0}, \lambda x. \lambda y. \square)$, using the recursor instead of pattern matching.) The next phrase “also suppose A or B ; we will derive a contradiction” indicates filling this hole by a function definition, introducing another unnamed hypothesis $z : A + B$, leading to the proof state:

$$f((x, y))(z) := \square : \mathbf{0}.$$

Now saying “there are two cases” indicates a case split, i.e. an application of the recursion principle for the coproduct $A + B$. If we write this using the recursor, it would be

$$f((x, y))(z) := \text{rec}_{A+B}(\mathbf{0}, \lambda a. \square, \lambda b. \square, z)$$

while if we write it using pattern matching, it would be

$$\begin{aligned} f((x, y))(\text{inl}(a)) &:= \square : \mathbf{0} \\ f((x, y))(\text{inr}(b)) &:= \square : \mathbf{0}. \end{aligned}$$

Note that in both cases we now have two “holes” of type $\mathbf{0}$ to fill in, corresponding to the two cases where we have to derive a contradiction. Finally, the conclusion of a contradiction from $a : A$ and $x : A \rightarrow \mathbf{0}$ is simply application of the function x to a , and similarly in the other case. (Note the convenient coincidence of the phrase “applying a function” with that of “applying a hypothesis” or theorem.) Thus our eventual definition is

$$\begin{aligned} f((x, y))(\text{inl}(a)) &:= x(a) \\ f((x, y))(\text{inr}(b)) &:= y(b). \end{aligned}$$

As an exercise, you should verify the converse tautology “If not (A or B), then (not A) and (not B)” by exhibiting an element of

$$((A + B) \rightarrow \mathbf{0}) \rightarrow (A \rightarrow \mathbf{0}) \times (B \rightarrow \mathbf{0}),$$

for any types A and B , using the rules we have just introduced.

However, not all classical tautologies hold under this interpretation. For example, the rule “If not (A and B), then (not A) or (not B)” is not valid: we cannot, in general, construct an element of the corresponding type

$$((A \times B) \rightarrow \mathbf{0}) \rightarrow (A \rightarrow \mathbf{0}) + (B \rightarrow \mathbf{0}).$$

This reflects the fact that the “natural” propositions-as-types logic of type theory is *constructive*. This means that it does not include certain classical principles, such as the law of excluded middle (LEM) or proof by contradiction, and others which depend on them, such as this instance of de Morgan’s law.

Philosophically, constructive logic is so-called because it confines itself to constructions that can be carried out *effectively*, which is to say those with a computational meaning. Without being too precise, this means there is some sort of algorithm specifying, step-by-step, how to build an object (and, as a special case, how to see that a theorem is true). This requires omission of LEM, since there is no *effective* procedure for deciding whether a proposition is true or false.

The constructivity of type-theoretic logic means it has an intrinsic computational meaning, which is of interest to computer scientists. It also means that type theory provides *axiomatic freedom*. For example, while by default there is no construction witnessing LEM, the logic is still compatible with the existence of one (see §3.4). Thus, because type theory does not *deny* LEM, we may consistently add it as an assumption, and work conventionally without restriction. In this respect, type theory enriches, rather than constrains, conventional mathematical practice.

We encourage the reader who is unfamiliar with constructive logic to work through some more examples as a means of getting familiar with it. See Exercises 1.12 and 1.13 for some suggestions.

So far we have discussed only propositional logic. Now we consider *predicate logic*, where in addition to logical connectives like “and” and “or” we have quantifiers “there exists” and “for all”. In this case, types play a dual role: they serve as propositions and also as types in the conventional sense, i.e., domains we quantify over. A predicate over a type A is represented as a family $P : A \rightarrow \mathcal{U}$, assigning to every element $a : A$ a type $P(a)$ corresponding to the proposition that P holds for a . We now extend the above translation with an explanation of the quantifiers:

English	Type Theory
For all $x : A$, $P(x)$ holds	$\prod_{(x:A)} P(x)$
There exists $x : A$ such that $P(x)$	$\sum_{(x:A)} P(x)$

As before, we can show that tautologies of (constructive) predicate logic translate into inhabited types. For example, *If for all $x : A$, $P(x)$ and $Q(x)$ then (for all $x : A$, $P(x)$) and (for all $x : A$, $Q(x)$)* translates to

$$(\prod_{(x:A)} P(x) \times Q(x)) \rightarrow (\prod_{(x:A)} P(x)) \times (\prod_{(x:A)} Q(x)).$$

An informal proof of this tautology might go as follows:

Suppose for all x , $P(x)$ and $Q(x)$. First, we suppose given x and prove $P(x)$. By assumption, we have $P(x)$ and $Q(x)$, and hence we have $P(x)$. Second, we suppose given x and prove $Q(x)$. Again by assumption, we have $P(x)$ and $Q(x)$, and hence we have $Q(x)$.

The first sentence begins defining an implication as a function, by introducing a witness for its hypothesis:

$$f(p) := \square : (\prod_{(x:A)} P(x)) \times (\prod_{(x:A)} Q(x)).$$

At this point there is an implicit use of the pairing constructor to produce an element of a product type, which is somewhat signposted in this example by the words “first” and “second”:

$$f(p) := \left(\square : \prod_{(x:A)} P(x), \square : \prod_{(x:A)} Q(x) \right).$$

The phrase “we suppose given x and prove $P(x)$ ” now indicates defining a *dependent* function in the usual way, introducing a variable for its input. Since this is inside a pairing constructor, it is natural to write it as a λ -abstraction:

$$f(p) : \equiv \left(\lambda x. (\square : P(x)) , \square : \prod_{(x:A)} Q(x) \right).$$

Now “we have $P(x)$ and $Q(x)$ ” invokes the hypothesis, obtaining $p(x) : P(x) \times Q(x)$, and “hence we have $P(x)$ ” implicitly applies the appropriate projection:

$$f(p) : \equiv \left(\lambda x. \text{pr}_1(p(x)) , \square : \prod_{(x:A)} Q(x) \right).$$

The next two sentences fill the other hole in the obvious way:

$$f(p) : \equiv \left(\lambda x. \text{pr}_1(p(x)) , \lambda x. \text{pr}_2(p(x)) \right).$$

Of course, the English proofs we have been using as examples are much more verbose than those that mathematicians usually use in practice; they are more like the sort of language one uses in an “introduction to proofs” class. The practicing mathematician has learned to fill in the gaps, so in practice we can omit plenty of details, and we will generally do so. The criterion of validity for proofs, however, is always that they can be translated back into the construction of an element of the corresponding type.

As a more concrete example, consider how to define inequalities of natural numbers. One natural definition is that $n \leq m$ if there exists a $k : \mathbb{N}$ such that $n + k = m$. (This uses again the identity types that we will introduce in the next section, but we will not need very much about them.) Under the propositions-as-types translation, this would yield:

$$(n \leq m) : \equiv \sum_{k:\mathbb{N}} (n + k = m).$$

The reader is invited to prove the familiar properties of \leq from this definition. For strict inequality, there are a couple of natural choices, such as

$$(n < m) : \equiv \sum_{k:\mathbb{N}} (n + \text{succ}(k) = m)$$

or

$$(n < m) : \equiv (n \leq m) \times \neg(n = m).$$

The former is more natural in constructive mathematics, but in this case it is actually equivalent to the latter, since \mathbb{N} has “decidable equality” (see §3.4 and Theorem 7.2.6).

There is also another interpretation of the type $\sum_{(x:A)} P(x)$. Since an inhabitant of it is an element $x : A$ together with a witness that $P(x)$ holds, instead of regarding $\sum_{(x:A)} P(x)$ as the proposition “there exists an $x : A$ such that $P(x)$ ”, we can regard it as “the type of all elements $x : A$ such that $P(x)$ ”, i.e. as a “subtype” of A .

We will return to this interpretation in §3.5. For now, we note that it allows us to incorporate axioms into the definition of types as mathematical structures which we discussed in §1.6. For example, suppose we want to define a **semigroup** to be a type A equipped with a binary operation $m : A \rightarrow A \rightarrow A$ (that is, a magma) and such that for all $x, y, z : A$ we have $m(x, m(y, z)) = m(m(x, y), z)$. This latter proposition is represented by the type

$$\prod_{x,y,z:A} m(x, m(y, z)) = m(m(x, y), z),$$

so the type of semigroups is

$$\text{Semigroup} := \sum_{(A:\mathcal{U})} \sum_{(m:A \rightarrow A \rightarrow A)} \prod_{(x,y,z:A)} m(x, m(y, z)) = m(m(x, y), z),$$

i.e. the subtype of Magma consisting of the semigroups. From an inhabitant of Semigroup we can extract the carrier A , the operation m , and a witness of the axiom, by applying appropriate projections. We will return to this example in §2.14.

Note also that we can use the universes in type theory to represent “higher order logic”—that is, we can quantify over all propositions or over all predicates. For example, we can represent the proposition *for all properties* $P : A \rightarrow \mathcal{U}$, *if* $P(a)$ *then* $P(b)$ as

$$\prod_{P:A \rightarrow \mathcal{U}} P(a) \rightarrow P(b)$$

where $A : \mathcal{U}$ and $a, b : A$. However, *a priori* this proposition lives in a different, higher, universe than the propositions we are quantifying over; that is

$$\left(\prod_{P:A \rightarrow \mathcal{U}_i} P(a) \rightarrow P(b) \right) : \mathcal{U}_{i+1}.$$

We will return to this issue in §3.5.

We have described here a “proof-relevant” translation of propositions, where the proofs of disjunctions and existential statements carry some information. For instance, if we have an inhabitant of $A + B$, regarded as a witness of “ A or B ”, then we know whether it came from A or from B . Similarly, if we have an inhabitant of $\sum_{(x:A)} P(x)$, regarded as a witness of “there exists $x : A$ such that $P(x)$ ”, then we know what the element x is (it is the first projection of the given inhabitant).

As a consequence of the proof-relevant nature of this logic, we may have “ A if and only if B ” (which, recall, means $(A \rightarrow B) \times (B \rightarrow A)$), and yet the types A and B exhibit different behavior. For instance, it is easy to verify that “ \mathbb{N} if and only if $\mathbf{1}$ ”, and yet clearly \mathbb{N} and $\mathbf{1}$ differ in important ways. The statement “ \mathbb{N} if and only if $\mathbf{1}$ ” tells us only that *when regarded as a mere proposition*, the type \mathbb{N} represents the same proposition as $\mathbf{1}$ (in this case, the true proposition). We sometimes express “ A if and only if B ” by saying that A and B are **logically equivalent**. This is to be distinguished from the stronger notion of *equivalence of types* to be introduced in §2.4 and Chapter 4: although \mathbb{N} and $\mathbf{1}$ are logically equivalent, they are not equivalent types.

In Chapter 3 we will introduce a class of types called “mere propositions” for which equivalence and logical equivalence coincide. Using these types, we will introduce a modification to the above-described logic that is sometimes appropriate, in which the additional information contained in disjunctions and existentials is discarded.

Finally, we note that the propositions-as-types correspondence can be viewed in reverse, allowing us to regard any type A as a proposition, which we prove by exhibiting an element of A . Sometimes we will state this proposition as “ A is **inhabited**”. That is, when we say that A is inhabited, we mean that we have given a (particular) element of A , but that we are choosing not to give a name to that element. Similarly, to say that A is *not inhabited* is the same as to give an element of $\neg A$. In particular, the empty type $\mathbf{0}$ is obviously not inhabited, since $\neg\mathbf{0} \equiv (\mathbf{0} \rightarrow \mathbf{0})$ is inhabited by id_0 .¹⁰

¹⁰This should not be confused with the statement that type theory is consistent, which is the *meta-theoretic* claim that it is not possible to obtain an element of $\mathbf{0}$ by following the rules of type theory.

1.12 Identity types

While the previous constructions can be seen as generalizations of standard set theoretic constructions, our way of handling identity seems to be specific to type theory. According to the propositions-as-types conception, the *proposition* that two elements of the same type $a, b : A$ are equal must correspond to some *type*. Since this proposition depends on what a and b are, these **equality types** or **identity types** must be type families dependent on two copies of A .

We may write the family as $\text{Id}_A : A \rightarrow A \rightarrow \mathcal{U}$ (not to be mistaken for the identity function id_A), so that $\text{Id}_A(a, b)$ is the type representing the proposition of equality between a and b . Once we are familiar with propositions-as-types, however, it is convenient to also use the standard equality symbol for this; thus “ $a = b$ ” will also be a notation for the type $\text{Id}_A(a, b)$ corresponding to the proposition that a equals b . For clarity, we may also write “ $a =_A b$ ” to specify the type A . If we have an element of $a =_A b$, we may say that a and b are **equal**, or sometimes **propoositionally equal** if we want to emphasize that this is different from the judgmental equality $a \equiv b$ discussed in §1.1.

Just as we remarked in §1.11 that the propositions-as-types versions of “or” and “there exists” can include more information than just the fact that the proposition is true, nothing prevents the type $a = b$ from also including more information. Indeed, this is the cornerstone of the homotopical interpretation, where we regard witnesses of $a = b$ as *paths* or *equivalences* between a and b in the space A . Just as there can be more than one path between two points of a space, there can be more than one witness that two objects are equal. Put differently, we may regard $a = b$ as the type of *identifications* of a and b , and there may be many different ways in which a and b can be identified. We will return to the interpretation in Chapter 2; for now we focus on the basic rules for the identity type. Just like all the other types considered in this chapter, it will have rules for formation, introduction, elimination, and computation, which behave formally in exactly the same way.

The formation rule says that given a type $A : \mathcal{U}$ and two elements $a, b : A$, we can form the type $(a =_A b) : \mathcal{U}$ in the same universe. The basic way to construct an element of $a = b$ is to know that a and b are the same. Thus, the introduction rule is a dependent function

$$\text{refl} : \prod_{a:A} (a =_A a)$$

called **reflexivity**, which says that every element of A is equal to itself (in a specified way). We regard refl_a as being the constant path at the point a .

In particular, this means that if a and b are *judgmentally* equal, $a \equiv b$, then we also have an element $\text{refl}_a : a =_A b$. This is well-typed because $a \equiv b$ means that also the type $a =_A b$ is judgmentally equal to $a =_A a$, which is the type of refl_a .

The induction principle (i.e. the elimination rule) for the identity types is one of the most subtle parts of type theory, and crucial to the homotopy interpretation. We begin by considering an important consequence of it, the principle that “equals may be substituted for equals”, as expressed by the following:

Indiscernability of identicals: For every family

$$C : A \rightarrow \mathcal{U}$$

there is a function

$$f : \prod_{(x,y:A)} \prod_{(p:x=_Ay)} C(x) \rightarrow C(y)$$

such that

$$f(x, x, \text{refl}_x) : \equiv \text{id}_{C(x)}.$$

This says that every family of types C respects equality, in the sense that applying C to *equal* elements of A also results in a function between the resulting types. The displayed equality states that the function associated to reflexivity is the identity function (and we shall see that, in general, the function $f(x, y, p) : C(x) \rightarrow C(y)$ is always an equivalence of types).

Indiscernability of identicals can be regarded as a recursion principle for the identity type, analogous to those given for booleans and natural numbers above. Just as $\text{rec}_{\mathbb{N}}$ gives a specified map $\mathbb{N} \rightarrow C$ for any other type C of a certain sort, indiscernability of identicals gives a specified map from $x =_A y$ to certain other reflexive, binary relations on A , namely those of the form $C(x) \rightarrow C(y)$ for some unary predicate $C(x)$. We could also formulate a more general recursion principle with respect to reflexive relations of the more general form $C(x, y)$. However, in order to fully characterize the identity type, we must generalize this recursion principle to an induction principle, which not only considers maps out of $x =_A y$ but also families over it. Put differently, we consider not only allowing equals to be substituted for equals, but also taking into account the evidence p for the equality.

1.12.1 Path induction

The induction principle for the identity type is called **path induction**, in view of the homotopical interpretation to be explained in the introduction to Chapter 2. It can be seen as stating that the family of identity types is freely generated by the elements of the form $\text{refl}_x : x = x$.

Path induction: Given a family

$$C : \prod_{x,y:A} (x =_A y) \rightarrow \mathcal{U}$$

and a function

$$c : \prod_{x:A} C(x, x, \text{refl}_x),$$

there is a function

$$f : \prod_{(x,y:A)} \prod_{(p:x=_Ay)} C(x, y, p)$$

such that

$$f(x, x, \text{refl}_x) : \equiv c(x).$$

Note that just like the induction principles for products, coproducts, natural numbers, and so on, path induction allows us to define *specified* functions which exhibit appropriate computational behavior. Just as we have *the* function $f : \mathbb{N} \rightarrow C$ defined by recursion from $c_0 : C$ and $c_s : \mathbb{N} \rightarrow C \rightarrow C$, which moreover satisfies $f(0) \equiv c_0$ and $f(\text{succ}(n)) \equiv c_s(n, f(n))$, we have *the* function $f : \prod_{(x,y:A)} \prod_{(p:x=_Ay)} C(x, y, p)$ defined by path induction from $c : \prod_{(x:A)} C(x, x, \text{refl}_x)$, which moreover satisfies $f(x, x, \text{refl}_x) \equiv c(x)$.

To understand the meaning of this principle, consider first the simpler case when C does not depend on p . Then we have $C : A \rightarrow A \rightarrow \mathcal{U}$, which we may regard as a predicate depending on two elements of A . We are interested in knowing when the proposition $C(x, y)$ holds for some pair of elements $x, y : A$. In this case, the hypothesis of path induction says that we know $C(x, x)$ holds for all $x : A$, i.e. that if we evaluate C at the pair x, x , we get a true proposition — so C is a

reflexive relation. The conclusion then tells us that $C(x, y)$ holds whenever $x = y$. This is exactly the more general recursion principle for reflexive relations mentioned above.

The general, inductive form of the rule allows C to also depend on the witness $p : x = y$ to the identity between x and y . In the premise, we not only replace x, y by x, x , but also simultaneously replace p by reflexivity: to prove a property for all elements x, y and paths $p : x = y$ between them, it suffices to consider all the cases where the elements are x, x and the path is $\text{refl}_x : x = x$. If we were viewing types just as sets, it would be unclear what this buys us, but since there may be many different identifications $p : x = y$ between x and y , it makes sense to keep track of them in considering families over the type $x =_A y$. In Chapter 2 we will see that this is very important to the homotopy interpretation.

If we package up path induction into a single function, it takes the form:

$$\text{ind}_{=A} : \prod_{(C:\prod_{(x,y:A)}(x=_Ay)\rightarrow\mathcal{U})} \left(\prod_{(x:A)} C(x, x, \text{refl}_x) \right) \rightarrow \prod_{(x,y:A)} \prod_{(p:x=_Ay)} C(x, y, p)$$

with the equality

$$\text{ind}_{=A}(C, c, x, x, \text{refl}_x) := c(x).$$

The function $\text{ind}_{=A}$ is traditionally called J . We will show in Lemma 2.3.1 that indiscernability of identicals is an instance of path induction, and also give it a new name and notation.

Given a proof $p : a = b$, path induction requires us to replace *both* a and b with the same unknown element x ; thus in order to define an element of a family C , for all pairs of equal elements of A , it suffices to define it on the diagonal. In some proofs, however, it is simpler to make use of an equation $p : a = b$ by replacing all occurrences of b with a (or vice versa), because it is sometimes easier to do the remainder of the proof for the specific element a mentioned in the equality than for a general unknown x . This motivates a second induction principle for identity types, which says that the family of types $a =_A x$ is generated by the element $\text{refl}_a : a = a$. As we show below, this second principle is equivalent to the first; it is just sometimes a more convenient formulation.

Based path induction: Fix an element $a : A$, and suppose given a family

$$C : \prod_{x:A} (a =_A x) \rightarrow \mathcal{U}$$

and an element

$$c : C(a, \text{refl}_a).$$

Then we obtain a function

$$f : \prod_{(x:A)} \prod_{(p:a=x)} C(x, p)$$

such that

$$f(a, \text{refl}_a) := c.$$

Here, $C(x, p)$ is a family of types, where x is an element of A and p is an element of the identity type $a =_A x$, for fixed a in A . The based path induction principle says that to define an element of this family for all x and p , it suffices to consider just the case where x is a and p is $\text{refl}_a : a = a$.

Packaged as a function, based path induction becomes:

$$\text{ind}'_{=_A} : \prod_{(a:A)} \prod_{(C:\prod_{(x:A)}(a=_A x) \rightarrow \mathcal{U})} C(a, \text{refl}_a) \rightarrow \prod_{(x:A)} \prod_{(p:a=_A x)} C(x, p)$$

with the equality

$$\text{ind}'_{=_A}(a, C, c, a, \text{refl}_a) \equiv c.$$

Below, we show that path induction and based path induction are equivalent. Because of this, we will sometimes be sloppy and also refer to based path induction simply as “path induction”, relying on the reader to infer which principle is meant from the form of the proof.

Remark 1.12.1. Intuitively, the induction principle for the natural numbers expresses the fact that every natural number is either 0 or of the form $\text{succ}(n)$ for some natural number n , so that if we prove a property for these cases (with induction hypothesis in the second case), then we have proved it for all natural numbers. Similarly, the induction principle for $A + B$ expresses the fact that every element of $A + B$ is either of the form $\text{inl}(a)$ or $\text{inr}(b)$, and so on. Applying this same reading to path induction, we might say that path induction expresses the fact that every path is of the form refl_a , so that if we prove a property for reflexivity paths, then we have proved it for all paths.

However, this reading is quite confusing in the context of the homotopy interpretation of paths, where there may be many different ways in which two elements a and b can be identified, and therefore many different elements of the identity type! How can there be many different paths, but at the same time we have an induction principle asserting that the only path is reflexivity?

The key observation is that it is not the identity *type* that is inductively defined, but the identity *family*. In particular, path induction says that the *family* of types $(x =_A y)$, as x, y vary over all elements of A , is inductively defined by the elements of the form refl_x . This means that to give an element of any other family $C(x, y, p)$ dependent on a *generic* element (x, y, p) of the identity family, it suffices to consider the cases of the form (x, x, refl_x) . In the homotopy interpretation, this says that the type of triples (x, y, p) , where x and y are the endpoints of the path p (in other words, the Σ -type $\sum_{(x,y:A)}(x = y)$), is inductively generated by the constant loops at each point x . As we will see in Chapter 2, in homotopy theory the space corresponding to $\sum_{(x,y:A)}(x = y)$ is the *free path space* — the space of paths in A whose endpoints may vary — and it is in fact the case that any point of this space is homotopic to the constant loop at some point, since we can simply retract one of its endpoints along the given path. The analogous fact is also true in type theory: we can prove by path induction on $p : x = y$ that $(x, y, p) =_{\sum_{(x,y:A)}(x = y)} (x, x, \text{refl}_x)$.

Similarly, based path induction says that for a fixed $a : A$, the *family* of types $(a =_A y)$, as y varies over all elements of A , is inductively defined by the element refl_a . Thus, to give an element of any other family $C(y, p)$ dependent on a generic element (y, p) of this family, it suffices to consider the case (a, refl_a) . Homotopically, this expresses the fact that the space of paths starting at some chosen point (the *based path space* at that point, which type-theoretically is $\sum_{(y:A)}(a = y)$) is contractible to the constant loop on the chosen point. Again, the corresponding fact is also true in type theory: we can prove by based path induction on $p : a = y$ that $(y, p) =_{\sum_{(y:A)}(a = y)} (a, \text{refl}_a)$. Note also that according to the interpretation of Σ -types as subtypes mentioned in §1.11, the type $\sum_{(y:A)}(a = y)$ can be regarded as “the type of all elements of A which are equal to a ”, a type-theoretic version of the “singleton subset” $\{a\}$.

Neither path induction nor based path induction provides a way to give an element of a family $C(p)$ where p has *two fixed endpoints* a and b . In particular, for a family $C : (a =_A a) \rightarrow \mathcal{U}$ dependent on a loop, we *cannot* apply path induction and consider only the case for $C(\text{refl}_a)$, and consequently, we cannot prove that all loops are reflexivity. Thus, inductively defining the identity family does not prohibit non-reflexivity paths in specific instances of the identity type. In other words, a path $p : x = x$ may be not equal to reflexivity as an element of $(x = x)$, but the pair (x, p) will nevertheless be equal to the pair (x, refl_x) as elements of $\sum_{(y:A)} (x = y)$.

As a topological example, consider a loop in the punctured disc $\{ (x, y) \mid 0 < x^2 + y^2 < 2 \}$ which starts at $(1, 0)$ and goes around the hole at $(0, 0)$ once before returning back to $(1, 0)$. If we hold both endpoints fixed at $(1, 0)$, this loop cannot be deformed into a constant path while staying within the punctured disc, just as a rope looped around a pole cannot be pulled in if we keep hold of both ends. However, the loop can be contracted back to a constant if we allow one endpoint to vary, just as we can always gather in a rope if we only hold onto one end.

1.12.2 Equivalence of path induction and based path induction

The two induction principles for the identity type introduced above are equivalent. It is easy to see that path induction follows from the based path induction principle. Indeed, let us assume the premises of path induction:

$$\begin{aligned} C &: \prod_{x,y:A} (x =_A y) \rightarrow \mathcal{U}, \\ c &: \prod_{x:A} C(x, x, \text{refl}_x). \end{aligned}$$

Now, given an element $x : A$, we can instantiate both of the above, obtaining

$$\begin{aligned} C' &: \prod_{y:A} (x =_A y) \rightarrow \mathcal{U}, \\ C' &\equiv C(x), \\ c' &: C'(x, \text{refl}_x), \\ c' &\equiv c(x). \end{aligned}$$

Clearly, C' and c' match the premises of based path induction and hence we can construct

$$g : \prod_{(y:A)} \prod_{(p:x=y)} C'(y, p)$$

with the defining equality

$$g(x, \text{refl}_x) \equiv c'.$$

Now we observe that g 's codomain is equal to $C(x, y, p)$. Thus, discharging our assumption $x : A$, we can derive a function

$$f : \prod_{(x,y:A)} \prod_{(p:x=_A y)} C(x, y, p)$$

with the required judgmental equality $f(x, x, \text{refl}_x) \equiv g(x, \text{refl}_x) \equiv c' \equiv c(x)$.

Another proof of this fact is to observe that any such f can be obtained as an instance of $\text{ind}'_{=_A}$ so it suffices to define $\text{ind}_{=_A}$ in terms of $\text{ind}'_{=_A}$ as

$$\text{ind}_{=_A}(C, c, x, y, p) \equiv \text{ind}'_{=_A}(x, C(x), c(x), y, p).$$

The other direction is a bit trickier; it is not clear how we can use a particular instance of path induction to derive a particular instance of based path induction. What we can do instead is to construct one instance of path induction which shows all possible instantiations of based path induction at once. Define

$$\begin{aligned} D : \prod_{x,y:A} (x =_A y) &\rightarrow \mathcal{U}, \\ D(x,y,p) &\equiv \prod_{C:\prod_{(z:A)}(x=_Az)\rightarrow\mathcal{U}} C(x, \text{refl}_x) \rightarrow C(y, p). \end{aligned}$$

Then we can construct the function

$$\begin{aligned} d : \prod_{x:A} D(x, x, \text{refl}_x), \\ d &\equiv \lambda x. \lambda C. \lambda(c:C(x, \text{refl}_x)). c \end{aligned}$$

and hence using path induction obtain

$$f : \prod_{(x,y:A)} \prod_{(p:x=_Ay)} D(x, y, p)$$

with $f(x, x, \text{refl}_x) := d(x)$. Unfolding the definition of D , we can expand the type of f :

$$f : \prod_{(x,y:A)} \prod_{(p:x=_Ay)} \prod_{(C:\prod_{(z:A)}(x=_Az)\rightarrow\mathcal{U})} C(x, \text{refl}_x) \rightarrow C(y, p).$$

Now given $x : A$ and $p : a =_A x$, we can derive the conclusion of based path induction:

$$f(a, x, p, C, c) : C(x, p).$$

Notice that we also obtain the correct definitional equality.

Another proof is to observe that any use of based path induction is an instance of $\text{ind}'_{=A}$ and to define

$$\begin{aligned} \text{ind}'_{=A}(a, C, c, x, p) &\equiv \text{ind}_{=A} ((\lambda x, y. \lambda p. \prod_{(C:\prod_{(z:A)}(x=_Az)\rightarrow\mathcal{U})} C(x, \text{refl}_x) \rightarrow C(y, p)), \\ &\quad (\lambda x. \lambda C. \lambda d. d), a, x, p, C, c). \end{aligned}$$

Note that the construction given above uses universes. That is, if we want to model $\text{ind}'_{=A}$ with $C : \prod_{(x:A)}(a =_A x) \rightarrow \mathcal{U}_i$, we need to use $\text{ind}_{=A}$ with

$$D : \prod_{x,y:A} (x =_A y) \rightarrow \mathcal{U}_{i+1}$$

since D quantifies over all C of the given type. While this is compatible with our definition of universes, it is also possible to derive $\text{ind}'_{=A}$ without using universes: we can show that $\text{ind}_{=A}$ entails Lemmas 2.3.1 and 3.11.8, and that these two principles imply $\text{ind}'_{=A}$ directly. We leave the details to the reader as Exercise 1.7.

We can use either of the foregoing formulations of identity types to establish that equality is an equivalence relation, that every function preserves equality and that every family respects equality. We leave the details to the next chapter, where this will be derived and explained in the context of homotopy type theory.

Remark 1.12.2. We emphasize that despite having some unfamiliar features, propositional equality is *the* equality of mathematics in homotopy type theory. This distinction does not belong to judgmental equality, which is rather a metatheoretic feature of the rules of type theory. For instance, the associativity of addition for natural numbers proven in §1.9 is a *propositional* equality, not a judgmental one. The same is true of the commutative law (Exercise 1.16). Even the very simple commutativity $n + 1 = 1 + n$ is not a judgmental equality for a generic n (though it is judgmental for any specific n , e.g. $3 + 1 \equiv 1 + 3$, since both are judgmentally equal to 4 by the computation rules defining $+$). We can only prove such facts by using the identity type, since we can only apply the induction principle for \mathbb{N} with a type as output (not a judgment).

1.12.3 Disequality

Finally, let us also say something about **disequality**, which is negation of equality:¹¹

$$(x \neq_A y) : \equiv \neg(x =_A y).$$

If $x \neq y$, we say that x and y are **unequal** or **not equal**. Just like negation, disequality plays a less important role here than it does in classical mathematics. For example, we cannot prove that two things are equal by proving that they are not unequal: that would be an application of the classical law of double negation, see §3.4.

Sometimes it is useful to phrase disequality in a positive way. For example, in ?? we shall prove that a real number x has an inverse if, and only if, its distance from 0 is positive, which is a stronger requirement than $x \neq 0$.

Notes

The type theory presented here is a version of Martin-Löf’s intuitionistic type theory [ML98, ML75, ML82, ML84], which itself is based on and influenced by the foundational work of Brouwer [Bee85], Heyting [Hey66], Scott [Sco70], de Bruijn [dB73], Howard [How80], Tait [Tai67, Tai68], and Lawvere [Law06]. Three principal variants of Martin-Löf’s type theory underlie the NUPRL [CAB⁺86], COQ [Coq12], and AGDA [Nor07] computer implementations of type theory. The theory given here differs from these formulations in a number of respects, some of which are critical to the homotopy interpretation, while others are technical conveniences or involve concepts that have not yet been studied in the homotopical setting.

Most significantly, the type theory described here is derived from the *intensional* version of Martin-Löf’s type theory [ML75], rather than the *extensional* version [ML82]. Whereas the extensional theory makes no distinction between judgmental and propositional equality, the intensional theory regards judgmental equality as purely definitional, and admits a much broader, proof-relevant interpretation of the identity type that is central to the homotopy interpretation. From the homotopical perspective, extensional type theory confines itself to homotopically discrete sets (see §3.1), whereas the intensional theory admits types with higher-dimensional structure. The NUPRL system [CAB⁺86] is extensional, whereas both COQ [Coq12] and AGDA [Nor07] are intensional. Among intensional type theories, there are a number of variants that differ in the structure of identity proofs. The most liberal interpretation, on which we rely here, admits a *proof-relevant* interpretation of equality, whereas more restricted variants impose restrictions

¹¹We use “inequality” to refer to $<$ and \leq . Also, note that this is negation of the *propositional* identity type. Of course, it makes no sense to negate judgmental equality \equiv , because judgments are not subject to logical operations.

such as *uniqueness of identity proofs* (*UIP*) [Str93], stating that any two proofs of equality are judgmentally equal, and *Axiom K* [Str93], stating that the only proof of equality is reflexivity (up to judgmental equality). These additional requirements may be selectively imposed in the Coq and AGDA systems.

Another point of variation among intensional theories is the strength of judgmental equality, particularly as regards objects of function type. Here we include the uniqueness principle (η -conversion) $f \equiv \lambda x. f(x)$, as a principle of judgmental equality. This principle is used, for example, in §4.9, to show that univalence implies propositional function extensionality. Uniqueness principles are sometimes considered for other types. For instance, the uniqueness principle for the cartesian product $A \times B$ would be a judgmental version of the propositional equality $\text{uniq}_{A \times B}$ which we constructed in §1.5, saying that $u \equiv (\text{pr}_1(u), \text{pr}_2(u))$. This and the corresponding version for dependent pairs would be reasonable choices (which we did not make), but we cannot include all such rules, because the corresponding uniqueness principle for identity types would trivialize all the higher homotopical structure. So we are *forced* to leave it out, and the question then becomes where to draw the line. With regards to inductive types, we discuss these points further in §5.5.

It is important for our purposes that (propositional) equality of functions is taken to be *extensional* (in a different sense than that used above!). This is not a consequence of the rules in this chapter; it will be expressed by Axiom 2.9.3. This decision is significant for our purposes, because it specifies that equality of functions is as expected in mathematics. Although we include Axiom 2.9.3 as an axiom, it may be derived from the univalence axiom and the uniqueness principle for functions (see §4.9), as well as from the existence of an interval type (see Lemma 6.3.2).

Regarding inductive types such as products, Σ -types, coproducts, natural numbers, and so on (see Chapter 5), there are additional choices regarding the formulation of induction and recursion. We have taken *induction principles* as basic and *pattern matching* as derived from them, but one may also do the other; see ???. Usually in the latter case one allows also *deep* pattern matching; see [Coq92b]. There are several reasons for our choice. One reason is that induction principles are what we obtain naturally in categorical semantics. Another is that specifying the allowable kinds of (deep) pattern matching is quite tricky; for instance, AGDA’s pattern matching can prove Axiom K by default, although a flag `--without-K` prevents this [CDP14]. Finally, deep pattern matching is not well-understood for *higher* inductive types (see Chapter 6). Therefore, we will only use pattern matches such as those described in §1.10, which are directly equivalent to the application of an induction principle.

Unlike the type theory of Coq, we do not include a primitive type of propositions. Instead, as discussed in §1.11, we embrace the *propositions-as-types* (*PAT*) principle, identifying propositions with types. This was suggested originally by de Bruijn [dB73], Howard [How80], Tait [Tai68], and Martin-Löf [ML98]. (Our decision is explained more fully in §§3.2 and 3.3.)

We do, however, include a full cumulative hierarchy of universes, so that the type formation and equality judgments become instances of the membership and equality judgments for a universe. As a convenience, we regard objects of a universe as types, rather than as codes for types; in the terminology of [ML84], this means we use “Russell-style universes” rather than “Tarski-style universes”. An alternative would be to use Tarski-style universes, with an explicit coercion function required to make an element $A : \mathcal{U}$ of a universe into a type $\text{El}(A)$, and just say that the coercion is omitted when working informally.

We also treat the universe hierarchy as cumulative, in that every type in \mathcal{U}_i is also in \mathcal{U}_j for each $j \geq i$. There are different ways to implement cumulativity formally: the simplest is just to include a rule that if $A : \mathcal{U}_i$ then $A : \mathcal{U}_j$. However, this has the annoying consequence that

for a type family $B : A \rightarrow \mathcal{U}_i$ we cannot conclude $B : A \rightarrow \mathcal{U}_j$, although we can conclude $\lambda a. B(a) : A \rightarrow \mathcal{U}_j$. A more sophisticated approach that solves this problem is to introduce a judgmental subtyping relation $<$ generated by $\mathcal{U}_i < \mathcal{U}_j$, but this makes the type theory more complicated to study. Another alternative would be to include an explicit coercion function $\uparrow : \mathcal{U}_i \rightarrow \mathcal{U}_j$, which could be omitted when working informally.

It is also not necessary that the universes be indexed by natural numbers and linearly ordered. For some purposes, it is more appropriate to assume only that every universe is an element of some larger universe, together with a “directedness” property that any two universes are jointly contained in some larger one. There are many other possible variations, such as including a universe “ \mathcal{U}_ω ” that contains all \mathcal{U}_i (or even higher “large cardinal” type universes), or by internalizing the hierarchy into a type family $\lambda i. \mathcal{U}_i$. The latter is in fact done in AGDA.

The path induction principle for identity types was formulated by Martin-Löf [ML98]. The based path induction rule in the setting of Martin-Löf type theory is due to Paulin-Mohring [PM93]; it can be seen as an intensional generalization of the concept of “pointwise functionality” for hypothetical judgments from NUPRL [CAB⁺86, Section 8.1]. The fact that Martin-Löf’s rule implies Paulin-Mohring’s was proved by Streicher using Axiom K (see §7.2), by Altenkirch and Goguen as in §1.12, and finally by Hofmann without universes (as in Exercise 1.7); see [Str93, §1.3 and Addendum].

Exercises

Exercise 1.1. Given functions $f : A \rightarrow B$ and $g : B \rightarrow C$, define their **composite** $g \circ f : A \rightarrow C$. Show that we have $h \circ (g \circ f) \equiv (h \circ g) \circ f$.

Exercise 1.2. Derive the recursion principle for products $\text{rec}_{A \times B}$ using only the projections, and verify that the definitional equalities are valid. Do the same for Σ -types.

Exercise 1.3. Derive the induction principle for products $\text{ind}_{A \times B}$, using only the projections and the propositional uniqueness principle $\text{uniq}_{A \times B}$. Verify that the definitional equalities are valid. Generalize $\text{uniq}_{A \times B}$ to Σ -types, and do the same for Σ -types. (*This requires concepts from Chapter 2.*)

Exercise 1.4. Assuming as given only the *iterator* for natural numbers

$$\text{iter} : \prod_{C:\mathcal{U}} C \rightarrow (C \rightarrow C) \rightarrow \mathbb{N} \rightarrow C$$

with the defining equations

$$\begin{aligned} \text{iter}(C, c_0, c_s, 0) &\equiv c_0, \\ \text{iter}(C, c_0, c_s, \text{succ}(n)) &\equiv c_s(\text{iter}(C, c_0, c_s, n)), \end{aligned}$$

derive a function having the type of the recursor $\text{rec}_{\mathbb{N}}$. Show that the defining equations of the recursor hold propositionally for this function, using the induction principle for \mathbb{N} .

Exercise 1.5. Show that if we define $A + B := \sum_{(x:2)} \text{rec}_2(\mathcal{U}, A, B, x)$, then we can give a definition of ind_{A+B} for which the definitional equalities stated in §1.7 hold.

Exercise 1.6. Show that if we define $A \times B := \prod_{(x:2)} \text{rec}_2(\mathcal{U}, A, B, x)$, then we can give a definition of $\text{ind}_{A \times B}$ for which the definitional equalities stated in §1.5 hold propositionally (i.e. using equality types). (*This requires the function extensionality axiom, which is introduced in §2.9.*)

Exercise 1.7. Give an alternative derivation of $\text{ind}'_{=A}$ from $\text{ind}_{=A}$ which avoids the use of universes. (*This is easiest using concepts from later chapters.*)

Exercise 1.8. Define multiplication and exponentiation using $\text{rec}_{\mathbb{N}}$. Verify that $(\mathbb{N}, +, 0, \times, 1)$ is a semiring using only $\text{ind}_{\mathbb{N}}$. You will probably also need to use symmetry and transitivity of equality, Lemmas 2.1.1 and 2.1.2.

Exercise 1.9. Define the type family $\text{Fin} : \mathbb{N} \rightarrow \mathcal{U}$ mentioned at the end of §1.3, and the dependent function $\text{fmax} : \prod_{(n:\mathbb{N})} \text{Fin}(n+1)$ mentioned in §1.4.

Exercise 1.10. Show that the Ackermann function $\text{ack} : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathbb{N}$ is definable using only $\text{rec}_{\mathbb{N}}$ satisfying the following equations:

$$\begin{aligned}\text{ack}(0, n) &\equiv \text{succ}(n), \\ \text{ack}(\text{succ}(m), 0) &\equiv \text{ack}(m, 1), \\ \text{ack}(\text{succ}(m), \text{succ}(n)) &\equiv \text{ack}(m, \text{ack}(\text{succ}(m), n)).\end{aligned}$$

Exercise 1.11. Show that for any type A , we have $\neg\neg\neg A \rightarrow \neg A$.

Exercise 1.12. Using the propositions as types interpretation, derive the following tautologies.

- (i) If A , then (if B then A).
- (ii) If A , then not (not A).
- (iii) If (not A or not B), then not (A and B).

Exercise 1.13. Using propositions-as-types, derive the double negation of the principle of excluded middle, i.e. prove *not (not (P or not P))*.

Exercise 1.14. Why do the induction principles for identity types not allow us to construct a function $f : \prod_{(x:A)} \prod_{(p:x=x)} (p = \text{refl}_x)$ with the defining equation

$$f(x, \text{refl}_x) : \equiv \text{refl}_{\text{refl}_x} \quad ?$$

Exercise 1.15. Show that indiscernability of identicals follows from path induction.

Exercise 1.16. Show that addition of natural numbers is commutative: $\prod_{(i,j:\mathbb{N})} (i + j = j + i)$.

Chapter 2

Homotopy type theory

The central new idea in homotopy type theory is that types can be regarded as spaces in homotopy theory, or higher-dimensional groupoids in category theory.

We begin with a brief summary of the connection between homotopy theory and higher-dimensional category theory. In classical homotopy theory, a space X is a set of points equipped with a topology, and a path between points x and y is represented by a continuous map $p : [0, 1] \rightarrow X$, where $p(0) = x$ and $p(1) = y$. This function can be thought of as giving a point in X at each “moment in time”. For many purposes, strict equality of paths (meaning, pointwise equal functions) is too fine a notion. For example, one can define operations of path concatenation (if p is a path from x to y and q is a path from y to z , then the concatenation $p \cdot q$ is a path from x to z) and inverses (p^{-1} is a path from y to x). However, there are natural equations between these operations that do not hold for strict equality: for example, the path $p \cdot p^{-1}$ (which walks from x to y , and then back along the same route, as time goes from 0 to 1) is not strictly equal to the identity path (which stays still at x at all times).

The remedy is to consider a coarser notion of equality of paths called *homotopy*. A homotopy between a pair of continuous maps $f : X_1 \rightarrow X_2$ and $g : X_1 \rightarrow X_2$ is a continuous map $H : X_1 \times [0, 1] \rightarrow X_2$ satisfying $H(x, 0) = f(x)$ and $H(x, 1) = g(x)$. In the specific case of paths p and q from x to y , a homotopy is a continuous map $H : [0, 1] \times [0, 1] \rightarrow X$ such that $H(s, 0) = p(s)$ and $H(s, 1) = q(s)$ for all $s \in [0, 1]$. In this case we require also that $H(0, t) = x$ and $H(1, t) = y$ for all $t \in [0, 1]$, so that for each t the function $H(-, t)$ is again a path from x to y ; a homotopy of this sort is said to be *endpoint-preserving* or *rel endpoints*. In simple cases, we can think of the image of the square $[0, 1] \times [0, 1]$ under H as “filling the space” between p and q , although for general X this doesn’t really make sense; it is better to think of H as a continuous deformation of p into q that doesn’t move the endpoints. Since $[0, 1] \times [0, 1]$ is 2-dimensional, we also speak of H as a 2-dimensional *path between paths*.

For example, because $p \cdot p^{-1}$ walks out and back along the same route, you know that you can continuously shrink $p \cdot p^{-1}$ down to the identity path—it won’t, for example, get snagged around a hole in the space. Homotopy is an equivalence relation, and operations such as concatenation, inverses, etc., respect it. Moreover, the homotopy equivalence classes of loops at some point x_0 (where two loops p and q are equated when there is a *based* homotopy between them, which is a homotopy H as above that additionally satisfies $H(0, t) = H(1, t) = x_0$ for all t) form a group called the *fundamental group*. This group is an *algebraic invariant* of a space, which can be used to investigate whether two spaces are *homotopy equivalent* (there are continuous maps back and forth whose composites are homotopic to the identity), because equivalent spaces have

isomorphic fundamental groups.

Because homotopies are themselves a kind of 2-dimensional path, there is a natural notion of 3-dimensional *homotopy between homotopies*, and then *homotopy between homotopies between homotopies*, and so on. This infinite tower of points, paths, homotopies, homotopies between homotopies, ..., equipped with algebraic operations such as the fundamental group, is an instance of an algebraic structure called a (weak) ∞ -groupoid. An ∞ -groupoid consists of a collection of objects, and then a collection of *morphisms* between objects, and then *morphisms between morphisms*, and so on, equipped with some complex algebraic structure; a morphism at level k is called a **k -morphism**. Morphisms at each level have identity, composition, and inverse operations, which are weak in the sense that they satisfy the groupoid laws (associativity of composition, identity is a unit for composition, inverses cancel) only up to morphisms at the next level, and this weakness gives rise to further structure. For example, because associativity of composition of morphisms $p \cdot (q \cdot r) = (p \cdot q) \cdot r$ is itself a higher-dimensional morphism, one needs an additional operation relating various proofs of associativity: the various ways to reassociate $p \cdot (q \cdot (r \cdot s))$ into $((p \cdot q) \cdot r) \cdot s$ give rise to Mac Lane's pentagon. Weakness also creates non-trivial interactions between levels.

Every topological space X has a *fundamental ∞ -groupoid* whose k -morphisms are the k -dimensional paths in X . The weakness of the ∞ -groupoid corresponds directly to the fact that paths form a group only up to homotopy, with the $(k+1)$ -paths serving as the homotopies between the k -paths. Moreover, the view of a space as an ∞ -groupoid preserves enough aspects of the space to do homotopy theory: the fundamental ∞ -groupoid construction is adjoint to the geometric realization of an ∞ -groupoid as a space, and this adjunction preserves homotopy theory (this is called the *homotopy hypothesis/theorem*, because whether it is a hypothesis or theorem depends on how you define ∞ -groupoid). For example, you can easily define the fundamental group of an ∞ -groupoid, and if you calculate the fundamental group of the fundamental ∞ -groupoid of a space, it will agree with the classical definition of fundamental group of that space. Because of this correspondence, homotopy theory and higher-dimensional category theory are intimately related.

Now, in homotopy type theory each type can be seen to have the structure of an ∞ -groupoid. Recall that for any type A , and any $x, y : A$, we have an identity type $x =_A y$, also written $\text{Id}_A(x, y)$ or just $x = y$. Logically, we may think of elements of $x = y$ as evidence that x and y are equal, or as identifications of x with y . Furthermore, type theory (unlike, say, first-order logic) allows us to consider such elements of $x =_A y$ also as individuals which may be the subjects of further propositions. Therefore, we can *iterate* the identity type: we can form the type $p =_{(x=_Ay)} q$ of identifications between identifications p, q , and the type $r =_{(p=(x=_Ay)q)} s$, and so on. The structure of this tower of identity types corresponds precisely to that of the continuous paths and (higher) homotopies between them in a space, or an ∞ -groupoid.

Thus, we will frequently refer to an element $p : x =_A y$ as a **path** from x to y ; we call x its **start point** and y its **end point**. Two paths $p, q : x =_A y$ with the same start and end point are said to be **parallel**, in which case an element $r : p =_{(x=_Ay)} q$ can be thought of as a homotopy, or a morphism between morphisms; we will often refer to it as a **2-path** or a **2-dimensional path**. Similarly, $r =_{(p=(x=_Ay)q)} s$ is the type of **3-dimensional paths** between two parallel 2-dimensional paths, and so on. If the type A is “set-like”, such as \mathbb{N} , these iterated identity types will be uninteresting (see §3.1), but in the general case they can model non-trivial homotopy types.

An important difference between homotopy type theory and classical homotopy theory is that homotopy type theory provides a *synthetic* description of spaces, in the following sense.

Synthetic geometry is geometry in the style of Euclid [EucBC]: one starts from some basic notions (points and lines), constructions (a line connecting any two points), and axioms (all right angles are equal), and deduces consequences logically. This is in contrast with analytic geometry, where notions such as points and lines are represented concretely using cartesian coordinates in \mathbb{R}^n —lines are sets of points—and the basic constructions and axioms are derived from this representation. While classical homotopy theory is analytic (spaces and paths are made of points), homotopy type theory is synthetic: points, paths, and paths between paths are basic, indivisible, primitive notions.

Moreover, one of the amazing things about homotopy type theory is that all of the basic constructions and axioms—all of the higher groupoid structure—arises automatically from the induction principle for identity types. Recall from §1.12 that this says that if

- for every $x, y : A$ and every $p : x =_A y$ we have a type $D(x, y, p)$, and
- for every $a : A$ we have an element $d(a) : D(a, a, \text{refl}_a)$,

then

- there exists an element $\text{ind}_{=A}(D, d, x, y, p) : D(x, y, p)$ for every two elements $x, y : A$ and $p : x =_A y$, such that $\text{ind}_{=A}(D, d, a, a, \text{refl}_a) \equiv d(a)$.

In other words, given dependent functions

$$\begin{aligned} D &: \prod_{x, y : A} (x = y) \rightarrow \mathcal{U} \\ d &: \prod_{a : A} D(a, a, \text{refl}_a) \end{aligned}$$

there is a dependent function

$$\text{ind}_{=A}(D, d) : \prod_{(x, y : A)} \prod_{(p : x = y)} D(x, y, p)$$

such that

$$\text{ind}_{=A}(D, d, a, a, \text{refl}_a) \equiv d(a) \tag{2.0.1}$$

for every $a : A$. Usually, every time we apply this induction rule we will either not care about the specific function being defined, or we will immediately give it a different name.

Informally, the induction principle for identity types says that if we want to construct an object (or prove a statement) which depends on an inhabitant $p : x =_A y$ of an identity type, then it suffices to perform the construction (or the proof) in the special case when x and y are the same (judgmentally) and p is the reflexivity element $\text{refl}_x : x = x$ (judgmentally). When writing informally, we may express this with a phrase such as “by induction, it suffices to assume...”. This reduction to the “reflexivity case” is analogous to the reduction to the “base case” and “inductive step” in an ordinary proof by induction on the natural numbers, and also to the “left case” and “right case” in a proof by case analysis on a disjoint union or disjunction.

The “conversion rule” (2.0.1) is less familiar in the context of proof by induction on natural numbers, but there is an analogous notion in the related concept of definition by recursion. If a sequence $(a_n)_{n \in \mathbb{N}}$ is defined by giving a_0 and specifying a_{n+1} in terms of a_n , then in fact the 0th term of the resulting sequence is the given one, and the given recurrence relation relating a_{n+1} to a_n holds for the resulting sequence. (This may seem so obvious as to not be worth saying, but if we view a definition by recursion as an algorithm for calculating values of a sequence, then

it is precisely the process of executing that algorithm.) The rule (2.0.1) is analogous: it says that if we define an object $f(p)$ for all $p : x = y$ by specifying what the value should be when p is $\text{refl}_x : x = x$, then the value we specified is in fact the value of $f(\text{refl}_x)$.

This induction principle endows each type with the structure of an ∞ -groupoid, and each function between two types with the structure of an ∞ -functor between two such groupoids. This is interesting from a mathematical point of view, because it gives a new way to work with ∞ -groupoids. It is interesting from a type-theoretic point of view, because it reveals new operations that are associated with each type and function. In the remainder of this chapter, we begin to explore this structure.

2.1 Types are higher groupoids

We now derive from the induction principle the beginnings of the structure of a higher groupoid. We begin with symmetry of equality, which, in topological language, means that “paths can be reversed”.

Lemma 2.1.1. *For every type A and every $x, y : A$ there is a function*

$$(x = y) \rightarrow (y = x)$$

*denoted $p \mapsto p^{-1}$, such that $\text{refl}_x^{-1} \equiv \text{refl}_x$ for each $x : A$. We call p^{-1} the **inverse** of p .*

Since this is our first time stating something as a “Lemma” or “Theorem”, let us pause to consider what that means. Recall that propositions (statements susceptible to proof) are identified with types, whereas lemmas and theorems (statements that have been proven) are identified with *inhabited* types. Thus, the statement of a lemma or theorem should be translated into a type, as in §1.11, and its proof translated into an inhabitant of that type. According to the interpretation of the universal quantifier “for every”, the type corresponding to Lemma 2.1.1 is

$$\prod_{(A:\mathcal{U})} \prod_{(x,y:A)} (x = y) \rightarrow (y = x).$$

The proof of Lemma 2.1.1 will consist of constructing an element of this type, i.e. deriving the judgment $f : \prod_{(A:\mathcal{U})} \prod_{(x,y:A)} (x = y) \rightarrow (y = x)$ for some f . We then introduce the notation $(-)^{-1}$ for this element f , in which the arguments A , x , and y are omitted and inferred from context. (As remarked in §1.1, the secondary statement “ $\text{refl}_x^{-1} \equiv \text{refl}_x$ for each $x : A$ ” should be regarded as a separate judgment.)

First proof. Assume given $A : \mathcal{U}$, and let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family defined by $D(x, y, p) := (y = x)$. In other words, D is a function assigning to any $x, y : A$ and $p : x = y$ a type, namely the type $y = x$. Then we have an element

$$d := \lambda x. \text{refl}_x : \prod_{x:A} D(x, x, \text{refl}_x).$$

Thus, the induction principle for identity types gives us an element $\text{ind}_{=A}(D, d, x, y, p) : (y = x)$ for each $p : (x = y)$. We can now define the desired function $(-)^{-1}$ to be $\lambda p. \text{ind}_{=A}(D, d, x, y, p)$, i.e. we set $p^{-1} := \text{ind}_{=A}(D, d, x, y, p)$. The conversion rule (2.0.1) gives $\text{refl}_x^{-1} \equiv \text{refl}_x$, as required. \square

We have written out this proof in a very formal style, which may be helpful while the induction rule on identity types is unfamiliar. To be even more formal, we could say that Lemma 2.1.1 and its proof together consist of the judgment

$$\lambda A. \lambda x. \lambda y. \lambda p. \text{ind}_{=A}((\lambda x. \lambda y. \lambda p. (y = x)), (\lambda x. \text{refl}_x), x, y, p) : \prod_{(A:\mathcal{U})} \prod_{(x,y:A)} (x = y) \rightarrow (y = x)$$

(along with an additional equality judgment). However, eventually we prefer to use more natural language, such as in the following equivalent proof.

Second proof. We want to construct, for each $x, y : A$ and $p : x = y$, an element $p^{-1} : y = x$. By induction, it suffices to do this in the case when y is x and p is refl_x . But in this case, the type $x = y$ of p and the type $y = x$ in which we are trying to construct p^{-1} are both simply $x = x$. Thus, in the “reflexivity case”, we can define refl_x^{-1} to be simply refl_x . The general case then follows by the induction principle, and the conversion rule $\text{refl}_x^{-1} \equiv \text{refl}_x$ is precisely the proof in the reflexivity case that we gave. \square

We will write out the next few proofs in both styles, to help the reader become accustomed to the latter one. Next we prove the transitivity of equality, or equivalently we “concatenate paths”.

Lemma 2.1.2. *For every type A and every $x, y, z : A$ there is a function*

$$(x = y) \rightarrow (y = z) \rightarrow (x = z)$$

written $p \mapsto q \mapsto p \cdot q$, such that $\text{refl}_x \cdot \text{refl}_x \equiv \text{refl}_x$ for any $x : A$. We call $p \cdot q$ the **concatenation** or **composite** of p and q .

First proof. Let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family

$$D(x, y, p) := \prod_{(z:A)} \prod_{(q:y=z)} (x = z).$$

Note that $D(x, x, \text{refl}_x) \equiv \prod_{(z:A)} \prod_{(q:x=z)} (x = z)$. Thus, in order to apply the induction principle for identity types to this D , we need a function of type

$$\prod_{x:A} D(x, x, \text{refl}_x) \tag{2.1.3}$$

which is to say, of type

$$\prod_{(x,z:A)} \prod_{(q:x=z)} (x = z).$$

Now let $E : \prod_{(x,z:A)} \prod_{(q:x=z)} \mathcal{U}$ be the type family $E(x, z, q) := (x = z)$. Note that $E(x, x, \text{refl}_x) \equiv (x = x)$. Thus, we have the function

$$e(x) := \text{refl}_x : E(x, x, \text{refl}_x).$$

By the induction principle for identity types applied to E , we obtain a function

$$d : \prod_{(x,z:A)} \prod_{(q:x=z)} E(x, z, q).$$

But $E(x, z, q) \equiv (x = z)$, so the type of d is (2.1.3). Thus, we can use this function d and apply the induction principle for identity types to D , to obtain our desired function of type

$$\prod_{x,y,z:A} (y = z) \rightarrow (x = y) \rightarrow (x = z).$$

The conversion rules for the two induction principles give us $\text{refl}_x \cdot \text{refl}_x \equiv \text{refl}_x$ for any $x : A$. \square

Second proof. We want to construct, for every $x, y, z : A$ and every $p : x = y$ and $q : y = z$, an element of $x = z$. By induction on p , it suffices to assume that y is x and p is refl_x . In this case, the type $y = z$ of q is $x = z$. Now by induction on q , it suffices to assume also that z is x and q is refl_x . But in this case, $x = z$ is $x = x$, and we have $\text{refl}_x : (x = x)$. \square

The reader may well feel that we have given an overly convoluted proof of this lemma. In fact, we could stop after the induction on p , since at that point what we want to produce is an equality $x = z$, and we already have such an equality, namely q . Why do we go on to do another induction on q ?

The answer is that, as described in the introduction, we are doing *proof-relevant* mathematics. When we prove a lemma, we are defining an inhabitant of some type, and it can matter what *specific* element we defined in the course of the proof, not merely the type inhabited by that element (that is, the *statement* of the lemma). Lemma 2.1.2 has three obvious proofs: we could do induction over p , induction over q , or induction over both of them. If we proved it three different ways, we would have three different elements of the same type. It's not hard to show that these three elements are equal (see Exercise 2.1), but as they are not *definitionally* equal, there can still be reasons to prefer one over another.

In the case of Lemma 2.1.2, the difference hinges on the computation rule. If we proved the lemma using a single induction over p , then we would end up with a computation rule of the form $\text{refl}_y \cdot q \equiv q$. If we proved it with a single induction over q , we would have instead $p \cdot \text{refl}_y \equiv p$, while proving it with a double induction (as we did) gives only $\text{refl}_x \cdot \text{refl}_x \equiv \text{refl}_x$.

The asymmetrical computation rules can sometimes be convenient when doing formalized mathematics, as they allow the computer to simplify more things automatically. However, in informal mathematics, and arguably even in the formalized case, it can be confusing to have a concatenation operation which behaves asymmetrically and to have to remember which side is the “special” one. Treating both sides symmetrically makes for more robust proofs; this is why we have given the proof that we did. (However, this is admittedly a stylistic choice.)

The table below summarizes the “equality”, “homotopical”, and “higher-groupoid” points of view on what we have done so far.

Equality	Homotopy	∞ -Groupoid
reflexivity	constant path	identity morphism
symmetry	inversion of paths	inverse morphism
transitivity	concatenation of paths	composition of morphisms

In practice, transitivity is often applied to prove an equality by a chain of intermediate steps. We will use the common notation for this such as $a = b = c = d$. If the intermediate expressions are long, or we want to specify the witness of each equality, we may write

$$\begin{aligned} a &= b && (\text{by } p) \\ &= c && (\text{by } q) \\ &= d && (\text{by } r). \end{aligned}$$

In either case, the notation indicates construction of the element $(p \cdot q) \cdot r : (a = d)$. (We choose left-associativity for concreteness, although in view of Lemma 2.1.4(iv) below it makes little dif-

ference.) If it should happen that b and c , say, are judgmentally equal, then we may write

$$\begin{aligned} a &= b && \text{(by } p\text{)} \\ &\equiv c \\ &= d && \text{(by } r\text{)} \end{aligned}$$

to indicate construction of $p \bullet r : (a = d)$. We also follow common mathematical practice in not requiring the justifications in this notation (“by p ” and “by r ”) to supply the exact witness needed; instead we allow them to simply mention the most important (or least obvious) ingredient in constructing that witness. For instance, if “Lemma A” states that for all x and y we have $f(x) = g(y)$, then we may write “by Lemma A” as a justification for the step $f(a) = g(b)$, trusting the reader to deduce that we apply Lemma A with $x := a$ and $y := b$. We may also omit a justification entirely if we trust the reader to be able to guess it.

Now, because of proof-relevance, we can’t stop after proving “symmetry” and “transitivity” of equality: we need to know that these *operations* on equalities are well-behaved. (This issue is invisible in set theory, where symmetry and transitivity are mere *properties* of equality, rather than structure on paths.) From the homotopy-theoretic point of view, concatenation and inversion are just the “first level” of higher groupoid structure — we also need coherence laws on these operations, and analogous operations at higher dimensions. For instance, we need to know that concatenation is *associative*, and that inversion provides *inverses* with respect to concatenation.

Lemma 2.1.4. *Suppose $A : \mathcal{U}$, that $x, y, z, w : A$ and that $p : x = y$ and $q : y = z$ and $r : z = w$. We have the following:*

- (i) $p = p \bullet \text{refl}_y$ and $p = \text{refl}_x \bullet p$.
- (ii) $p^{-1} \bullet p = \text{refl}_y$ and $p \bullet p^{-1} = \text{refl}_x$.
- (iii) $(p^{-1})^{-1} = p$.
- (iv) $p \bullet (q \bullet r) = (p \bullet q) \bullet r$.

Note, in particular, that (i)–(iv) are themselves propositional equalities, living in the identity types of identity types, such as $p =_{x=y} q$ for $p, q : x = y$. Topologically, they are *paths of paths*, i.e. homotopies. It is a familiar fact in topology that when we concatenate a path p with the reversed path p^{-1} , we don’t literally obtain a constant path (which corresponds to the equality refl in type theory) — instead we have a homotopy, or higher path, from $p \bullet p^{-1}$ to the constant path.

Proof of Lemma 2.1.4. All the proofs use the induction principle for equalities.

(i) *First proof:* let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family given by

$$D(x, y, p) := (p = p \bullet \text{refl}_y).$$

Then $D(x, x, \text{refl}_x)$ is $\text{refl}_x = \text{refl}_x \bullet \text{refl}_x$. Since $\text{refl}_x \bullet \text{refl}_x \equiv \text{refl}_x$, it follows that $D(x, x, \text{refl}_x) \equiv (\text{refl}_x = \text{refl}_x)$. Thus, there is a function

$$d := \lambda x. \text{refl}_{\text{refl}_x} : \prod_{x:A} D(x, x, \text{refl}_x).$$

Now the induction principle for identity types gives an element $\text{ind}_{=A}(D, d, x, y, p) : (p = p \bullet \text{refl}_y)$ for each $p : x = y$. The other equality is proven similarly.

Second proof: by induction on p , it suffices to assume that y is x and that p is refl_x . But in this case, we have $\text{refl}_x \bullet \text{refl}_x \equiv \text{refl}_x$.

(ii) *First proof:* let $D : \prod_{(x,y:A)}(x = y) \rightarrow \mathcal{U}$ be the type family given by

$$D(x, y, p) := (p^{-1} \bullet p = \text{refl}_y).$$

Then $D(x, x, \text{refl}_x)$ is $\text{refl}_x^{-1} \bullet \text{refl}_x = \text{refl}_x$. Since $\text{refl}_x^{-1} \equiv \text{refl}_x$ and $\text{refl}_x \bullet \text{refl}_x \equiv \text{refl}_x$, we get that $D(x, x, \text{refl}_x) \equiv (\text{refl}_x = \text{refl}_x)$. Hence we find the function

$$d := \lambda x. \text{refl}_{\text{refl}_x} : \prod_{x:A} D(x, x, \text{refl}_x).$$

Now path induction gives an element $\text{ind}_{=A}(D, d, x, y, p) : p^{-1} \bullet p = \text{refl}_y$ for each $p : x = y$ in A . The other equality is similar.

Second proof: by induction, it suffices to assume p is refl_x . But in this case, we have $p^{-1} \bullet p \equiv \text{refl}_x^{-1} \bullet \text{refl}_x \equiv \text{refl}_x$.

(iii) *First proof:* let $D : \prod_{(x,y:A)}(x = y) \rightarrow \mathcal{U}$ be the type family given by

$$D(x, y, p) := (p^{-1}^{-1} = p).$$

Then $D(x, x, \text{refl}_x)$ is the type $(\text{refl}_x^{-1}^{-1} = \text{refl}_x)$. But since $\text{refl}_x^{-1} \equiv \text{refl}_x$ for each $x : A$, we have $\text{refl}_x^{-1}^{-1} \equiv \text{refl}_x^{-1} \equiv \text{refl}_x$, and thus $D(x, x, \text{refl}_x) \equiv (\text{refl}_x = \text{refl}_x)$. Hence we find the function

$$d := \lambda x. \text{refl}_{\text{refl}_x} : \prod_{x:A} D(x, x, \text{refl}_x).$$

Now path induction gives an element $\text{ind}_{=A}(D, d, x, y, p) : p^{-1}^{-1} = p$ for each $p : x = y$.

Second proof: by induction, it suffices to assume p is refl_x . But in this case, we have $p^{-1}^{-1} \equiv \text{refl}_x^{-1}^{-1} \equiv \text{refl}_x$.

(iv) *First proof:* let $D_1 : \prod_{(x,y:A)}(x = y) \rightarrow \mathcal{U}$ be the type family given by

$$D_1(x, y, p) := \prod_{(z,w:A)} \prod_{(q:y=z)} \prod_{(r:z=w)} (p \bullet (q \bullet r) = (p \bullet q) \bullet r).$$

Then $D_1(x, x, \text{refl}_x)$ is

$$\prod_{(z,w:A)} \prod_{(q:x=z)} \prod_{(r:z=w)} (\text{refl}_x \bullet (q \bullet r) = (\text{refl}_x \bullet q) \bullet r).$$

To construct an element of this type, let $D_2 : \prod_{(x,z:A)}(x = z) \rightarrow \mathcal{U}$ be the type family

$$D_2(x, z, q) := \prod_{(w:A)} \prod_{(r:z=w)} (\text{refl}_x \bullet (q \bullet r) = (\text{refl}_x \bullet q) \bullet r).$$

Then $D_2(x, x, \text{refl}_x)$ is

$$\prod_{(w:A)} \prod_{(r:x=w)} (\text{refl}_x \bullet (\text{refl}_x \bullet r) = (\text{refl}_x \bullet \text{refl}_x) \bullet r).$$

To construct an element of this type, let $D_3 : \prod_{(x,w:A)}(x = w) \rightarrow \mathcal{U}$ be the type family

$$D_3(x, w, r) := (\text{refl}_x \bullet (\text{refl}_x \bullet r) = (\text{refl}_x \bullet \text{refl}_x) \bullet r).$$

Then $D_3(x, x, \text{refl}_x)$ is

$$(\text{refl}_x \bullet (\text{refl}_x \bullet \text{refl}_x)) = (\text{refl}_x \bullet \text{refl}_x) \bullet \text{refl}_x$$

which is definitionally equal to the type $(\text{refl}_x = \text{refl}_x)$, and is therefore inhabited by $\text{refl}_{\text{refl}_x}$. Applying the path induction rule three times, therefore, we obtain an element of the overall desired type.

Second proof: by induction, it suffices to assume p, q , and r are all refl_x . But in this case, we have

$$\begin{aligned} p \bullet (q \bullet r) &\equiv \text{refl}_x \bullet (\text{refl}_x \bullet \text{refl}_x) \\ &\equiv \text{refl}_x \\ &\equiv (\text{refl}_x \bullet \text{refl}_x) \bullet \text{refl}_x \\ &\equiv (p \bullet q) \bullet r. \end{aligned}$$

Thus, we have $\text{refl}_{\text{refl}_x}$ inhabiting this type. \square

Remark 2.1.5. There are other ways to define these higher paths. For instance, in Lemma 2.1.4(iv) we might do induction only over one or two paths rather than all three. Each possibility will produce a *definitionally* different proof, but they will all be equal to each other. Such an equality between any two particular proofs can, again, be proven by induction, reducing all the paths in question to reflexivities and then observing that both proofs reduce themselves to reflexivities.

In view of Lemma 2.1.4(iv), we will often write $p \bullet q \bullet r$ for $(p \bullet q) \bullet r$, and similarly $p \bullet q \bullet r \bullet s$ for $((p \bullet q) \bullet r) \bullet s$ and so on. We choose left-associativity for definiteness, but it makes no real difference. We generally trust the reader to insert instances of Lemma 2.1.4(iv) to reassociate such expressions as necessary.

We are still not really done with the higher groupoid structure: the paths (i)–(iv) must also satisfy their own higher coherence laws, which are themselves higher paths, and so on “all the way up to infinity” (this can be made precise using e.g. the notion of a globular operad). However, for most purposes it is unnecessary to make the whole infinite-dimensional structure explicit. One of the nice things about homotopy type theory is that all of this structure can be *proven* starting from only the inductive property of identity types, so we can make explicit as much or as little of it as we need.

In particular, in this book we will not need any of the complicated combinatorics involved in making precise notions such as “coherent structure at all higher levels”. In addition to ordinary paths, we will use paths of paths (i.e. elements of a type $p =_{x=Ay} q$), which as remarked previously we call *2-paths* or *2-dimensional paths*, and perhaps occasionally paths of paths of paths (i.e. elements of a type $r =_{p=x=Ayq} s$), which we call *3-paths* or *3-dimensional paths*. It is possible to define a general notion of *n-dimensional path* (see Exercise 2.4), but we will not need it.

We will, however, use one particularly important and simple case of higher paths, which is when the start and end points are the same. In set theory, the proposition $a = a$ is entirely uninteresting, but in homotopy theory, paths from a point to itself are called *loops* and carry lots of interesting higher structure. Thus, given a type A with a point $a : A$, we define its **loop space** $\Omega(A, a)$ to be the type $a =_A a$. We may sometimes write simply ΩA if the point a is understood from context.

Since any two elements of ΩA are paths with the same start and end points, they can be concatenated; thus we have an operation $\Omega A \times \Omega A \rightarrow \Omega A$. More generally, the higher groupoid structure of A gives ΩA the analogous structure of a “higher group”.

It can also be useful to consider the loop space of the loop space of A , which is the space of 2-dimensional loops on the identity loop at a . This is written $\Omega^2(A, a)$ and represented in type theory by the type $\text{refl}_a =_{(a=_{AA}a)} \text{refl}_a$. While $\Omega^2(A, a)$, as a loop space, is again a “higher group”, it now also has some additional structure resulting from the fact that its elements are 2-dimensional loops between 1-dimensional loops.

Theorem 2.1.6 (Eckmann–Hilton). *The composition operation on the second loop space*

$$\Omega^2(A) \times \Omega^2(A) \rightarrow \Omega^2(A)$$

is commutative: $\alpha \cdot \beta = \beta \cdot \alpha$, for any $\alpha, \beta : \Omega^2(A)$.

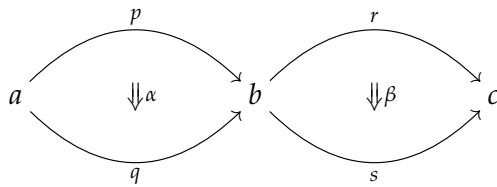
Proof. First, observe that the composition of 1-loops $\Omega A \times \Omega A \rightarrow \Omega A$ induces an operation

$$\star : \Omega^2(A) \times \Omega^2(A) \rightarrow \Omega^2(A)$$

as follows: consider elements $a, b, c : A$ and 1- and 2-paths,

$$\begin{array}{ll} p : a = b, & r : b = c \\ q : a = b, & s : b = c \\ \alpha : p = q, & \beta : r = s \end{array}$$

as depicted in the following diagram (with paths drawn as arrows).



Composing the upper and lower 1-paths, respectively, we get two paths $p \cdot r, q \cdot s : a = c$, and there is then a “horizontal composition”

$$\alpha \star \beta : p \cdot r = q \cdot s$$

between them, defined as follows. First, we define $\alpha \cdot_r r : p \cdot r = q \cdot r$ by path induction on r , so that

$$\alpha \cdot_r \text{refl}_b \equiv \text{ru}_p^{-1} \cdot \alpha \cdot \text{ru}_q$$

where $\text{ru}_p : p = p \cdot \text{refl}_b$ is the right unit law from Lemma 2.1.4(i). We could similarly define \cdot_r by induction on α , or on all paths in sight, resulting in different judgmental equalities, but for present purposes the definition by induction on r will make things simpler. Similarly, we define $q \cdot_l \beta : q \cdot r = q \cdot s$ by induction on q , so that

$$\text{refl}_b \cdot_l \beta \equiv \text{lu}_r^{-1} \cdot \beta \cdot \text{lu}_s$$

where lu_r denotes the left unit law. The operations \cdot_l and \cdot_r are called **whiskering**. Next, since $\alpha \cdot_r r$ and $q \cdot_l \beta$ are composable 2-paths, we can define the **horizontal composition** by:

$$\alpha \star \beta := (\alpha \cdot_r r) \cdot (q \cdot_l \beta).$$

Now suppose that $a \equiv b \equiv c$, so that all the 1-paths p, q, r , and s are elements of $\Omega(A, a)$, and assume moreover that $p \equiv q \equiv r \equiv s \equiv \text{refl}_a$, so that $\alpha : \text{refl}_a = \text{refl}_a$ and $\beta : \text{refl}_a = \text{refl}_a$ are composable in both orders. In that case, we have

$$\begin{aligned}\alpha * \beta &\equiv (\alpha \cdot_r \text{refl}_a) \cdot (\text{refl}_a \cdot_l \beta) \\ &= \text{ru}_{\text{refl}_a}^{-1} \cdot \alpha \cdot \text{ru}_{\text{refl}_a} \cdot \text{lu}_{\text{refl}_a}^{-1} \cdot \beta \cdot \text{lu}_{\text{refl}_a} \\ &\equiv \text{refl}_{\text{refl}_a}^{-1} \cdot \alpha \cdot \text{refl}_{\text{refl}_a} \cdot \text{refl}_{\text{refl}_a}^{-1} \cdot \beta \cdot \text{refl}_{\text{refl}_a} \\ &= \alpha \cdot \beta.\end{aligned}$$

(Recall that $\text{ru}_{\text{refl}_a} \equiv \text{lu}_{\text{refl}_a} \equiv \text{refl}_{\text{refl}_a}$, by the computation rule for path induction.) On the other hand, we can define another horizontal composition analogously by

$$\alpha \star' \beta := (p \cdot_l \beta) \cdot (\alpha \cdot_r s)$$

and we similarly learn that

$$\alpha \star' \beta = \beta \cdot \alpha.$$

But, in general, the two ways of defining horizontal composition agree, $\alpha * \beta = \alpha \star' \beta$, as we can see by induction on α and β and then on the two remaining 1-paths, to reduce everything to reflexivity. Thus we have

$$\alpha \cdot \beta = \alpha * \beta = \alpha \star' \beta = \beta \cdot \alpha. \quad \square$$

The foregoing fact, which is known as the *Eckmann–Hilton argument*, comes from classical homotopy theory, and indeed it is used in ?? below to show that the higher homotopy groups of a type are always abelian groups. The whiskering and horizontal composition operations defined in the proof are also a general part of the ∞ -groupoid structure of types. They satisfy their own laws (up to higher homotopy), such as

$$\alpha \cdot_r (p \cdot q) = (\alpha \cdot_r p) \cdot_r q$$

and so on. From now on, we trust the reader to apply path induction whenever needed to define further operations of this sort and verify their properties.

As this example suggests, the algebra of higher path types is much more intricate than just the groupoid-like structure at each level; the levels interact to give many further operations and laws, as in the study of iterated loop spaces in homotopy theory. Indeed, as in classical homotopy theory, we can make the following general definitions:

Definition 2.1.7. A **pointed type** (A, a) is a type $A : \mathcal{U}$ together with a point $a : A$, called its **basepoint**. We write $\mathcal{U}_\bullet := \sum_{(A:\mathcal{U})} A$ for the type of pointed types in the universe \mathcal{U} .

Definition 2.1.8. Given a pointed type (A, a) , we define the **loop space** of (A, a) to be the following pointed type:

$$\Omega(A, a) := ((a =_A a), \text{refl}_a).$$

An element of it will be called a **loop** at a . For $n : \mathbb{N}$, the **n -fold iterated loop space** $\Omega^n(A, a)$ of a pointed type (A, a) is defined recursively by:

$$\begin{aligned}\Omega^0(A, a) &:= (A, a) \\ \Omega^{n+1}(A, a) &:= \Omega^n(\Omega(A, a)).\end{aligned}$$

An element of it will be called an **n -loop** or an **n -dimensional loop** at a .

We will return to iterated loop spaces in Chapters 6 and 7 and ??.

2.2 Functions are functors

Now we wish to establish that functions $f : A \rightarrow B$ behave functorially on paths. In traditional type theory, this is equivalently the statement that functions respect equality. Topologically, this corresponds to saying that every function is “continuous”, i.e. preserves paths.

Lemma 2.2.1. *Suppose that $f : A \rightarrow B$ is a function. Then for any $x, y : A$ there is an operation*

$$\text{ap}_f : (x =_A y) \rightarrow (f(x) =_B f(y)).$$

Moreover, for each $x : A$ we have $\text{ap}_f(\text{refl}_x) \equiv \text{refl}_{f(x)}$.

The notation ap_f can be read either as the application of f to a path, or as the action on paths of f .

First proof. Let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family defined by

$$D(x, y, p) := (f(x) = f(y)).$$

Then we have

$$d := \lambda x. \text{refl}_{f(x)} : \prod_{x:A} D(x, x, \text{refl}_x).$$

By path induction, we obtain $\text{ap}_f : \prod_{(x,y:A)} (x = y) \rightarrow (f(x) = f(y))$. The computation rule implies $\text{ap}_f(\text{refl}_x) \equiv \text{refl}_{f(x)}$ for each $x : A$. \square

Second proof. To define $\text{ap}_f(p)$ for all $p : x = y$, it suffices, by induction, to assume p is refl_x . In this case, we may define $\text{ap}_f(p) := \text{refl}_{f(x)} : f(x) = f(x)$. \square

We will often write $\text{ap}_f(p)$ as simply $f(p)$. This is strictly speaking ambiguous, but generally no confusion arises. It matches the common convention in category theory of using the same symbol for the application of a functor to objects and to morphisms.

We note that ap behaves functorially, in all the ways that one might expect.

Lemma 2.2.2. *For functions $f : A \rightarrow B$ and $g : B \rightarrow C$ and paths $p : x =_A y$ and $q : y =_A z$, we have:*

- (i) $\text{ap}_f(p \bullet q) = \text{ap}_f(p) \bullet \text{ap}_f(q)$.
- (ii) $\text{ap}_f(p^{-1}) = \text{ap}_f(p)^{-1}$.
- (iii) $\text{ap}_g(\text{ap}_f(p)) = \text{ap}_{g \circ f}(p)$.
- (iv) $\text{ap}_{\text{id}_A}(p) = p$.

Proof. Left to the reader. \square

As was the case for the equalities in Lemma 2.1.4, those in Lemma 2.2.2 are themselves paths, which satisfy their own coherence laws (which can be proved in the same way), and so on.

2.3 Type families are fibrations

Since *dependently typed* functions are essential in type theory, we will also need a version of Lemma 2.2.1 for these. However, this is not quite so simple to state, because if $f : \prod_{(x:A)} B(x)$ and $p : x = y$, then $f(x) : B(x)$ and $f(y) : B(y)$ are elements of distinct types, so that *a priori* we cannot even ask whether they are equal. The missing ingredient is that p itself gives us a way to relate the types $B(x)$ and $B(y)$.

We have already seen this in section 1.12, where we called it “indiscernability of identicals”. We now introduce a different name and notation for it that we will use from now on.

Lemma 2.3.1 (Transport). *Suppose that P is a type family over A and that $p : x =_A y$. Then there is a function $p_* : P(x) \rightarrow P(y)$.*

First proof. Let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family defined by

$$D(x, y, p) := P(x) \rightarrow P(y).$$

Then we have the function

$$d := \lambda x. \text{id}_{P(x)} : \prod_{x:A} D(x, x, \text{refl}_x),$$

so that the induction principle gives us $\text{ind}_{=A}(D, d, x, y, p) : P(x) \rightarrow P(y)$ for $p : x = y$, which we define to be p_* . \square

Second proof. By induction, it suffices to assume p is refl_x . But in this case, we can take $(\text{refl}_x)_* : P(x) \rightarrow P(x)$ to be the identity function. \square

Sometimes, it is necessary to notate the type family P in which the transport operation happens. In this case, we may write

$$\text{transport}^P(p, -) : P(x) \rightarrow P(y).$$

Recall that a type family P over a type A can be seen as a property of elements of A , which holds at x in A if $P(x)$ is inhabited. Then the transportation lemma says that P respects equality, in the sense that if x is equal to y , then $P(x)$ holds if and only if $P(y)$ holds. In fact, we will see later on that if $x = y$ then actually $P(x)$ and $P(y)$ are equivalent.

Topologically, the transportation lemma can be viewed as a “path lifting” operation in a fibration. We think of a type family $P : A \rightarrow \mathcal{U}$ as a *fibration* with base space A , with $P(x)$ being the fiber over x , and with $\sum_{(x:A)} P(x)$ being the **total space** of the fibration, with first projection $\sum_{(x:A)} P(x) \rightarrow A$. The defining property of a fibration is that given a path $p : x = y$ in the base space A and a point $u : P(x)$ in the fiber over x , we may lift the path p to a path in the total space starting at u (and this lifting can be done continuously). The point $p_*(u)$ can be thought of as the other endpoint of this lifted path. We can also define the path itself in type theory:

Lemma 2.3.2 (Path lifting property). *Let $P : A \rightarrow \mathcal{U}$ be a type family over A and assume we have $u : P(x)$ for some $x : A$. Then for any $p : x = y$, we have*

$$\text{lift}(u, p) : (x, u) = (y, p_*(u))$$

in $\sum_{(x:A)} P(x)$, such that $\text{pr}_1(\text{lift}(u, p)) = p$.

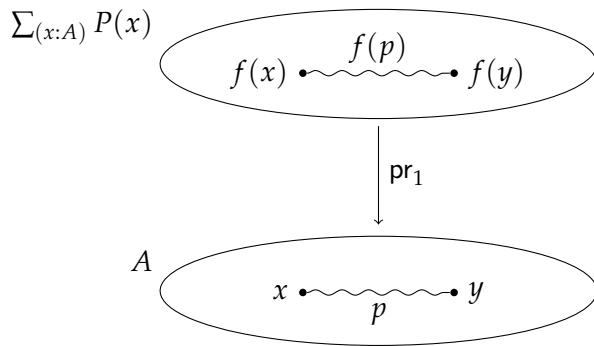
Proof. Left to the reader. We will prove a more general theorem in §2.7. \square

In classical homotopy theory, a fibration is defined as a map for which there *exist* liftings of paths; while in contrast, we have just shown that in type theory, every type family comes with a *specified* “path-lifting function”. This accords with the philosophy of constructive mathematics, according to which we cannot show that something exists except by exhibiting it. It also ensures automatically that the path liftings are chosen “continuously”, since as we have seen, all functions in type theory are “continuous”.

Remark 2.3.3. Although we may think of a type family $P : A \rightarrow \mathcal{U}$ as like a fibration, it is generally not a good idea to say things like “the fibration $P : A \rightarrow \mathcal{U}$ ”, since this sounds like we are talking about a fibration with base \mathcal{U} and total space A . To repeat, when a type family $P : A \rightarrow \mathcal{U}$ is regarded as a fibration, the base is A and the total space is $\sum_{(x:A)} P(x)$.

We may also occasionally use other topological terminology when speaking about type families. For instance, we may refer to a dependent function $f : \prod_{(x:A)} P(x)$ as a **section** of the fibration P , and we may say that something happens **fiberwise** if it happens for each $P(x)$. For instance, a section $f : \prod_{(x:A)} P(x)$ shows that P is “fiberwise inhabited”.

Now we can prove the dependent version of Lemma 2.2.1. The topological intuition is that given $f : \prod_{(x:A)} P(x)$ and a path $p : x =_A y$, we ought to be able to apply f to p and obtain a path in the total space of P which “lies over” p , as shown below.



We *can* obtain such a thing from Lemma 2.2.1. Given $f : \prod_{(x:A)} P(x)$, we can define a non-dependent function $f' : A \rightarrow \sum_{(x:A)} P(x)$ by setting $f'(x) := (x, f(x))$, and then consider $f'(p) : f'(x) = f'(y)$. Since $\text{pr}_1 \circ f' \equiv \text{id}_A$, by Lemma 2.2.2 we have $\text{pr}_1(f'(p)) = p$; thus $f'(p)$ does “lie over” p in this sense. However, it is not obvious from the *type* of $f'(p)$ that it lies over any specific path in A (in this case, p), which is sometimes important.

The solution is to use the transport lemma. By Lemma 2.3.2 we have a canonical path $\text{lift}(u, p)$ from (x, u) to $(y, p_*(u))$ which lies over p . Thus, any path from $u : P(x)$ to $v : P(y)$ lying over p should factor through $\text{lift}(u, p)$, essentially uniquely, by a path from $p_*(u)$ to v lying entirely in the fiber $P(y)$. Thus, up to equivalence, it makes sense to define “a path from u to v lying over $p : x = y$ ” to mean a path $p_*(u) = v$ in $P(y)$. And, indeed, we can show that dependent functions produce such paths.

Lemma 2.3.4 (Dependent map). *Suppose $f : \prod_{(x:A)} P(x)$; then we have a map*

$$\text{apd}_f : \prod_{p:x=y} (p_*(f(x)) =_{P(y)} f(y)).$$

First proof. Let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family defined by

$$D(x, y, p) := p_*(f(x)) = f(y).$$

Then $D(x, x, \text{refl}_x)$ is $(\text{refl}_x)_*(f(x)) = f(x)$. But since $(\text{refl}_x)_*(f(x)) \equiv f(x)$, we get that $D(x, x, \text{refl}_x) \equiv (f(x) = f(x))$. Thus, we find the function

$$d := \lambda x. \text{refl}_{f(x)} : \prod_{x:A} D(x, x, \text{refl}_x)$$

and now path induction gives us $\text{apd}_f(p) : p_*(f(x)) = f(y)$ for each $p : x = y$. \square

Second proof. By induction, it suffices to assume p is refl_x . But in this case, the desired equation is $(\text{refl}_x)_*(f(x)) = f(x)$, which holds judgmentally. \square

We will refer generally to paths which “lie over other paths” in this sense as *dependent paths*. They will play an increasingly important role starting in Chapter 6. In §2.5 we will see that for a few particular kinds of type families, there are equivalent ways to represent the notion of dependent paths that are sometimes more convenient.

Now recall from §1.4 that a non-dependently typed function $f : A \rightarrow B$ is just the special case of a dependently typed function $f : \prod_{(x:A)} P(x)$ when P is a constant type family, $P(x) := B$. In this case, apd_f and ap_f are closely related, because of the following lemma:

Lemma 2.3.5. *If $P : A \rightarrow \mathcal{U}$ is defined by $P(x) := B$ for a fixed $B : \mathcal{U}$, then for any $x, y : A$ and $p : x = y$ and $b : B$ we have a path*

$$\text{transportconst}_p^B(b) : \text{transport}^P(p, b) = b.$$

First proof. Fix a $b : B$, and let $D : \prod_{(x,y:A)} (x = y) \rightarrow \mathcal{U}$ be the type family defined by

$$D(x, y, p) := (\text{transport}^P(p, b) = b).$$

Then $D(x, x, \text{refl}_x)$ is $(\text{transport}^P(\text{refl}_x, b) = b)$, which is judgmentally equal to $(b = b)$ by the computation rule for transporting. Thus, we have the function

$$d := \lambda x. \text{refl}_b : \prod_{x:A} D(x, x, \text{refl}_x).$$

Now path induction gives us an element of $\prod_{(x,y:A)} \prod_{(p:x=y)} (\text{transport}^P(p, b) = b)$, as desired. \square

Second proof. By induction, it suffices to assume y is x and p is refl_x . But $\text{transport}^P(\text{refl}_x, b) \equiv b$, so in this case what we have to prove is $b = b$, and we have refl_b for this. \square

Thus, for any $x, y : A$ and $p : x = y$ and $f : A \rightarrow B$, by concatenating with $\text{transportconst}_p^B(f(x))$ and its inverse, respectively, we obtain functions

$$(f(x) = f(y)) \rightarrow (p_*(f(x)) = f(y)) \quad \text{and} \tag{2.3.6}$$

$$(p_*(f(x)) = f(y)) \rightarrow (f(x) = f(y)). \tag{2.3.7}$$

In fact, these functions are inverse equivalences (in the sense to be introduced in §2.4), and they relate $\text{ap}_f(p)$ to $\text{apd}_f(p)$.

Lemma 2.3.8. *For $f : A \rightarrow B$ and $p : x =_A y$, we have*

$$\text{apd}_f(p) = \text{transportconst}_p^B(f(x)) \bullet \text{ap}_f(p).$$

First proof. Let $D : \prod_{(x,y:A)}(x = y) \rightarrow \mathcal{U}$ be the type family defined by

$$D(x, y, p) := (\text{apd}_f(p) = \text{transportconst}_p^B(f(x)) \cdot \text{ap}_f(p)).$$

Thus, we have

$$D(x, x, \text{refl}_x) = (\text{apd}_f(\text{refl}_x) = \text{transportconst}_{\text{refl}_x}^B(f(x)) \cdot \text{ap}_f(\text{refl}_x)).$$

But by definition, all three paths appearing in this type are $\text{refl}_{f(x)}$, so we have

$$\text{refl}_{\text{refl}_{f(x)}} : D(x, x, \text{refl}_x).$$

Thus, path induction gives us an element of $\prod_{(x,y:A)} \prod_{(p:x=y)} D(x, y, p)$, which is what we wanted. \square

Second proof. By induction, it suffices to assume y is x and p is refl_x . In this case, what we have to prove is $\text{refl}_{f(x)} = \text{refl}_{f(x)} \cdot \text{refl}_{f(x)}$, which is true judgmentally. \square

Because the types of apd_f and ap_f are different, it is often clearer to use different notations for them.

At this point, we hope the reader is starting to get a feel for proofs by induction on identity types. From now on we stop giving both styles of proofs, allowing ourselves to use whatever is most clear and convenient (and often the second, more concise one). Here are a few other useful lemmas about transport; we leave it to the reader to give the proofs (in either style).

Lemma 2.3.9. *Given $P : A \rightarrow \mathcal{U}$ with $p : x =_A y$ and $q : y =_A z$ while $u : P(x)$, we have*

$$q_*(p_*(u)) = (p \cdot q)_*(u).$$

Lemma 2.3.10. *For a function $f : A \rightarrow B$ and a type family $P : B \rightarrow \mathcal{U}$, and any $p : x =_A y$ and $u : P(f(x))$, we have*

$$\text{transport}^{P \circ f}(p, u) = \text{transport}^P(\text{ap}_f(p), u).$$

Lemma 2.3.11. *For $P, Q : A \rightarrow \mathcal{U}$ and a family of functions $f : \prod_{(x:A)} P(x) \rightarrow Q(x)$, and any $p : x =_A y$ and $u : P(x)$, we have*

$$\text{transport}^Q(p, f_x(u)) = f_y(\text{transport}^P(p, u)).$$

2.4 Homotopies and equivalences

So far, we have seen how the identity type $x =_A y$ can be regarded as a type of *identifications*, *paths*, or *equivalences* between two elements x and y of a type A . Now we investigate the appropriate notions of “identification” or “sameness” between *functions* and between *types*. In §§2.9 and 2.10, we will see that homotopy type theory allows us to identify these with instances of the identity type, but before we can do that we need to understand them in their own right.

Traditionally, we regard two functions as the same if they take equal values on all inputs. Under the propositions-as-types interpretation, this suggests that two functions f and g (perhaps dependently typed) should be the same if the type $\prod_{(x:A)}(f(x) = g(x))$ is inhabited. Under the homotopical interpretation, this dependent function type consists of *continuous* paths or *functional* equivalences, and thus may be regarded as the type of *homotopies* or of *natural isomorphisms*. We will adopt the topological terminology for this.

Definition 2.4.1. Let $f, g : \prod_{(x:A)} P(x)$ be two sections of a type family $P : A \rightarrow \mathcal{U}$. A **homotopy** from f to g is a dependent function of type

$$(f \sim g) := \prod_{x:A} (f(x) = g(x)).$$

Note that a homotopy is not the same as an identification ($f = g$). However, in §2.9 we will introduce an axiom making homotopies and identifications “equivalent”.

The following proofs are left to the reader.

Lemma 2.4.2. *Homotopy is an equivalence relation on each dependent function type $\prod_{(x:A)} P(x)$. That is, we have elements of the types*

$$\begin{aligned} & \prod_{f:\prod_{(x:A)} P(x)} (f \sim f) \\ & \prod_{f,g:\prod_{(x:A)} P(x)} (f \sim g) \rightarrow (g \sim f) \\ & \prod_{f,g,h:\prod_{(x:A)} P(x)} (f \sim g) \rightarrow (g \sim h) \rightarrow (f \sim h). \end{aligned}$$

Just as functions in type theory are automatically “functors”, homotopies are automatically “natural transformations”. We will state and prove this only for non-dependent functions $f, g : A \rightarrow B$; in Exercise 2.18 we ask the reader to generalize it to dependent functions.

Recall that for $f : A \rightarrow B$ and $p : x =_A y$, we may write $f(p)$ to mean $\text{ap}_f(p)$.

Lemma 2.4.3. *Suppose $H : f \sim g$ is a homotopy between functions $f, g : A \rightarrow B$ and let $p : x =_A y$. Then we have*

$$H(x) \bullet g(p) = f(p) \bullet H(y).$$

We may also draw this as a commutative diagram:

$$\begin{array}{ccc} f(x) & \xlongequal{f(p)} & f(y) \\ H(x) \parallel & & \parallel H(y) \\ g(x) & \xlongequal{g(p)} & g(y) \end{array}$$

Proof. By induction, we may assume p is refl_x . Since ap_f and ap_g compute on reflexivity, in this case what we must show is

$$H(x) \bullet \text{refl}_{g(x)} = \text{refl}_{f(x)} \bullet H(x).$$

But this follows since both sides are equal to $H(x)$. \square

Corollary 2.4.4. *Let $H : f \sim \text{id}_A$ be a homotopy, with $f : A \rightarrow A$. Then for any $x : A$ we have*

$$H(f(x)) = f(H(x)).$$

Here $f(x)$ denotes the ordinary application of f to x , while $f(H(x))$ denotes $\text{ap}_f(H(x))$.

Proof. By naturality of H , the following diagram of paths commutes:

$$\begin{array}{ccc} ffx & \xrightarrow{f(Hx)} & fx \\ H(fx) \parallel & & \parallel Hx \\ fx & \xrightarrow{Hx} & x \end{array}$$

That is, $f(Hx) \cdot Hx = H(fx) \cdot Hx$. We can now whisker by $(Hx)^{-1}$ to cancel Hx , obtaining

$$f(Hx) = f(Hx) \cdot Hx \cdot (Hx)^{-1} = H(fx) \cdot Hx \cdot (Hx)^{-1} = H(fx)$$

as desired (with some associativity paths suppressed). \square

Of course, like the functoriality of functions (Lemma 2.2.2), the equality in Lemma 2.4.3 is a path which satisfies its own coherence laws, and so on.

Moving on to types, from a traditional perspective one may say that a function $f : A \rightarrow B$ is an *isomorphism* if there is a function $g : B \rightarrow A$ such that both composites $f \circ g$ and $g \circ f$ are pointwise equal to the identity, i.e. such that $f \circ g \sim \text{id}_B$ and $g \circ f \sim \text{id}_A$. A homotopical perspective suggests that this should be called a *homotopy equivalence*, and from a categorical one, it should be called an *equivalence of (higher) groupoids*. However, when doing proof-relevant mathematics, the corresponding type

$$\sum_{g:B \rightarrow A} ((f \circ g \sim \text{id}_B) \times (g \circ f \sim \text{id}_A)) \quad (2.4.5)$$

is poorly behaved. For instance, for a single function $f : A \rightarrow B$ there may be multiple unequal inhabitants of (2.4.5). (This is closely related to the observation in higher category theory that often one needs to consider *adjoint* equivalences rather than plain equivalences.) For this reason, we give (2.4.5) the following historically accurate, but slightly derogatory-sounding name instead.

Definition 2.4.6. For a function $f : A \rightarrow B$, a **quasi-inverse** of f is a triple (g, α, β) consisting of a function $g : B \rightarrow A$ and homotopies $\alpha : f \circ g \sim \text{id}_B$ and $\beta : g \circ f \sim \text{id}_A$.

Thus, (2.4.5) is *the type of quasi-inverses of f* ; we may denote it by $\text{qinv}(f)$.

Example 2.4.7. The identity function $\text{id}_A : A \rightarrow A$ has a quasi-inverse given by id_A itself, together with homotopies defined by $\alpha(y) := \text{refl}_y$ and $\beta(x) := \text{refl}_x$.

Example 2.4.8. For any $p : x =_A y$ and $z : A$, the functions

$$\begin{aligned} (p \bullet -) : (y =_A z) &\rightarrow (x =_A z) && \text{and} \\ (- \bullet p) : (z =_A x) &\rightarrow (z =_A y) \end{aligned}$$

have quasi-inverses given by $(p^{-1} \bullet -)$ and $(- \bullet p^{-1})$, respectively; see Exercise 2.6.

Example 2.4.9. For any $p : x =_A y$ and $P : A \rightarrow \mathcal{U}$, the function

$$\text{transport}^P(p, -) : P(x) \rightarrow P(y)$$

has a quasi-inverse given by $\text{transport}^P(p^{-1}, -)$; this follows from Lemma 2.3.9.

In general, we will only use the word *isomorphism* (and similar words such as *bijection*, and the associated notation $A \cong B$) in the special case when the types A and B “behave like sets” (see §3.1). In this case, the type (2.4.5) is unproblematic. We will reserve the word *equivalence* for an improved notion $\text{isequiv}(f)$ with the following properties:

- (i) For each $f : A \rightarrow B$ there is a function $\text{qinv}(f) \rightarrow \text{isequiv}(f)$.
- (ii) Similarly, for each f we have $\text{isequiv}(f) \rightarrow \text{qinv}(f)$; thus the two are logically equivalent (see §1.11).
- (iii) For any two inhabitants $e_1, e_2 : \text{isequiv}(f)$ we have $e_1 = e_2$.

In Chapter 4 we will see that there are many different definitions of $\text{isequiv}(f)$ which satisfy these three properties, but that all of them are equivalent. For now, to convince the reader that such things exist, we mention only the easiest such definition:

$$\text{isequiv}(f) := \left(\sum_{g:B \rightarrow A} (f \circ g \sim \text{id}_B) \right) \times \left(\sum_{h:B \rightarrow A} (h \circ f \sim \text{id}_A) \right). \quad (2.4.10)$$

We can show (i) and (ii) for this definition now. A function $\text{qinv}(f) \rightarrow \text{isequiv}(f)$ is easy to define by taking (g, α, β) to (g, α, g, β) . In the other direction, given (g, α, h, β) , let γ be the composite homotopy

$$g \xrightarrow{\beta} h \circ f \circ g \xrightarrow{\alpha} h,$$

meaning that $\gamma(x) := \beta(g(x))^{-1} \cdot h(\alpha(x))$. Now define $\beta' : g \circ f \sim \text{id}_A$ by $\beta'(x) := \gamma(f(x)) \cdot \beta(x)$. Then $(g, \alpha, \beta') : \text{qinv}(f)$.

Property (iii) for this definition is not too hard to prove either, but it requires identifying the identity types of cartesian products and dependent pair types, which we will discuss in §§2.6 and 2.7. Thus, we postpone it as well; see §4.3. At this point, the main thing to take away is that there is a well-behaved type which we can pronounce as “ f is an equivalence”, and that we can prove f to be an equivalence by exhibiting a quasi-inverse to it. In practice, this is the most common way to prove that a function is an equivalence.

In accord with the proof-relevant philosophy, an *equivalence* from A to B is defined to be a function $f : A \rightarrow B$ together with an inhabitant of $\text{isequiv}(f)$, i.e. a proof that it is an equivalence. We write $(A \simeq B)$ for the type of equivalences from A to B , i.e. the type

$$(A \simeq B) := \sum_{f:A \rightarrow B} \text{isequiv}(f). \quad (2.4.11)$$

Property (iii) above will ensure that if two equivalences are equal as functions (that is, the underlying elements of $A \rightarrow B$ are equal), then they are also equal as equivalences (see §2.7). Thus, we often abuse notation and blur the distinction between equivalences and their underlying functions. For instance, if we have a function $f : A \rightarrow B$ and we know that $e : \text{isequiv}(f)$, we may write $f : A \simeq B$, rather than (f, e) . Or conversely, if we have an equivalence $g : A \simeq B$, we may write $g(a)$ when given $a : A$, rather than $(\text{pr}_1 g)(a)$.

We conclude by observing:

Lemma 2.4.12. *Type equivalence is an equivalence relation on \mathcal{U} . More specifically:*

- (i) *For any A , the identity function id_A is an equivalence; hence $A \simeq A$.*
- (ii) *For any $f : A \simeq B$, we have an equivalence $f^{-1} : B \simeq A$.*
- (iii) *For any $f : A \simeq B$ and $g : B \simeq C$, we have $g \circ f : A \simeq C$.*

Proof. The identity function is clearly its own quasi-inverse; hence it is an equivalence.

If $f : A \rightarrow B$ is an equivalence, then it has a quasi-inverse, say $f^{-1} : B \rightarrow A$. Then f is also a quasi-inverse of f^{-1} , so f^{-1} is an equivalence $B \rightarrow A$.

Finally, given $f : A \simeq B$ and $g : B \simeq C$ with quasi-inverses f^{-1} and g^{-1} , say, then for any $a : A$ we have $f^{-1}g^{-1}gfa = f^{-1}fa = a$, and for any $c : C$ we have $gff^{-1}g^{-1}c = gg^{-1}c = c$. Thus $f^{-1} \circ g^{-1}$ is a quasi-inverse to $g \circ f$, hence the latter is an equivalence. \square

2.5 The higher groupoid structure of type formers

In Chapter 1, we introduced many ways to form new types: cartesian products, disjoint unions, dependent products, dependent sums, etc. In §§2.1–2.3, we saw that *all* types in homotopy type theory behave like spaces or higher groupoids. Our goal in the rest of the chapter is to make explicit how this higher structure behaves in the case of the particular types defined in Chapter 1.

It turns out that for many types A , the equality types $x =_A y$ can be characterized, up to equivalence, in terms of whatever data was used to construct A . For example, if A is a cartesian product $B \times C$, and $x \equiv (b, c)$ and $y \equiv (b', c')$, then we have an equivalence

$$((b, c) = (b', c')) \simeq ((b = b') \times (c = c')). \quad (2.5.1)$$

In more traditional language, two ordered pairs are equal just when their components are equal (but the equivalence (2.5.1) says rather more than this). The higher structure of the identity types can also be expressed in terms of these equivalences; for instance, concatenating two equalities between pairs corresponds to pairwise concatenation.

Similarly, when a type family $P : A \rightarrow \mathcal{U}$ is built up fiberwise using the type forming rules from Chapter 1, the operation $\text{transport}^P(p, -)$ can be characterized, up to homotopy, in terms of the corresponding operations on the data that went into P . For instance, if $P(x) \equiv B(x) \times C(x)$, then we have

$$\text{transport}^P(p, (b, c)) = (\text{transport}^B(p, b), \text{transport}^C(p, c)).$$

Finally, the type forming rules are also functorial, and if a function f is built from this functoriality, then the operations ap_f and apd_f can be computed based on the corresponding ones on the data going into f . For instance, if $g : B \rightarrow B'$ and $h : C \rightarrow C'$ and we define $f : B \times C \rightarrow B' \times C'$ by $f(b, c) := (g(b), h(c))$, then modulo the equivalence (2.5.1), we can identify ap_f with “ $(\text{ap}_g, \text{ap}_h)$ ”.

The next few sections (§§2.6–2.13) will be devoted to stating and proving theorems of this sort for all the basic type forming rules, with one section for each basic type former. Here we encounter a certain apparent deficiency in currently available type theories; as will become clear in later chapters, it would seem to be more convenient and intuitive if these characterizations of identity types, transport, and so on were *judgmental* equalities. However, in the theory presented in Chapter 1, the identity types are defined uniformly for all types by their induction principle, so we cannot “redefine” them to be different things at different types. Thus, the characterizations for particular types to be discussed in this chapter are, for the most part, *theorems* which we have to discover and prove, if possible.

Actually, the type theory of Chapter 1 is insufficient to prove the desired theorems for two of the type formers: Π -types and universes. For this reason, we are forced to introduce axioms into our type theory, in order to make those “theorems” true. Type-theoretically, an *axiom* (c.f. §1.1) is an “atomic” element that is declared to inhabit some specified type, without there being any rules governing its behavior other than those pertaining to the type it inhabits.

The axiom for Π -types (§2.9) is familiar to type theorists: it is called *function extensionality*, and states (roughly) that if two functions are homotopic in the sense of §2.4, then they are equal. The axiom for universes (§2.10), however, is a new contribution of homotopy type theory due to Voevodsky: it is called the *univalence axiom*, and states (roughly) that if two types are equivalent in the sense of §2.4, then they are equal. We have already remarked on this axiom in the introduction; it will play a very important role in this book.¹

It is important to note that not *all* identity types can be “determined” by induction over the construction of types. Counterexamples include most nontrivial higher inductive types (see Chapter 6 and ??). For instance, calculating the identity types of the types S^n (see §6.4) is equivalent to calculating the higher homotopy groups of spheres, a deep and important field of research in algebraic topology.

2.6 Cartesian product types

Given types A and B , consider the cartesian product type $A \times B$. For any elements $x, y : A \times B$ and a path $p : x =_{A \times B} y$, by functoriality we can extract paths $\text{pr}_1(p) : \text{pr}_1(x) =_A \text{pr}_1(y)$ and $\text{pr}_2(p) : \text{pr}_2(x) =_B \text{pr}_2(y)$. Thus, we have a function

$$(x =_{A \times B} y) \rightarrow (\text{pr}_1(x) =_A \text{pr}_1(y)) \times (\text{pr}_2(x) =_B \text{pr}_2(y)). \quad (2.6.1)$$

Theorem 2.6.2. *For any x and y , the function (2.6.1) is an equivalence.*

Read logically, this says that two pairs are equal just if they are equal componentwise. Read category-theoretically, this says that the morphisms in a product groupoid are pairs of morphisms. Read homotopy-theoretically, this says that the paths in a product space are pairs of paths.

Proof. We need a function in the other direction:

$$(\text{pr}_1(x) =_A \text{pr}_1(y)) \times (\text{pr}_2(x) =_B \text{pr}_2(y)) \rightarrow (x =_{A \times B} y). \quad (2.6.3)$$

By the induction rule for cartesian products, we may assume that x and y are both pairs, i.e. $x \equiv (a, b)$ and $y \equiv (a', b')$ for some $a, a' : A$ and $b, b' : B$. In this case, what we want is a function

$$(a =_A a') \times (b =_B b') \rightarrow ((a, b) =_{A \times B} (a', b')).$$

Now by induction for the cartesian product in its domain, we may assume given $p : a = a'$ and $q : b = b'$. And by two path inductions, we may assume that $a \equiv a'$ and $b \equiv b'$ and both p and q are reflexivity. But in this case, we have $(a, b) \equiv (a', b')$ and so we can take the output to also be reflexivity.

It remains to prove that (2.6.3) is quasi-inverse to (2.6.1). This is a simple sequence of inductions, but they have to be done in the right order.

In one direction, let us start with $r : x =_{A \times B} y$. We first do a path induction on r in order to assume that $x \equiv y$ and r is reflexivity. In this case, since ap_{pr_1} and ap_{pr_2} are defined by path induction, (2.6.1) takes $r \equiv \text{refl}_x$ to the pair $(\text{refl}_{\text{pr}_1 x}, \text{refl}_{\text{pr}_2 x})$. Now by induction on x , we may assume $x \equiv (a, b)$, so that this is $(\text{refl}_a, \text{refl}_b)$. Thus, (2.6.3) takes it by definition to $\text{refl}_{(a,b)}$, which (under our current assumptions) is r .

¹We have chosen to introduce these principles as axioms, but there are potentially other ways to formulate a type theory in which they hold. See the Notes to this chapter.

In the other direction, if we start with $s : (\text{pr}_1(x) =_A \text{pr}_1(y)) \times (\text{pr}_2(x) =_B \text{pr}_2(y))$, then we first do induction on x and y to assume that they are pairs (a, b) and (a', b') , and then induction on $s : (a =_A a') \times (b =_B b')$ to reduce it to a pair (p, q) where $p : a = a'$ and $q : b = b'$. Now by induction on p and q , we may assume they are reflexivities refl_a and refl_b , in which case (2.6.3) yields $\text{refl}_{(a,b)}$ and then (2.6.1) returns us to $(\text{refl}_a, \text{refl}_b) \equiv (p, q) \equiv s$. \square

In particular, we have shown that (2.6.1) has an inverse (2.6.3), which we may denote by

$$\text{pair}^= : (\text{pr}_1(x) = \text{pr}_1(y)) \times (\text{pr}_2(x) = \text{pr}_2(y)) \rightarrow (x = y).$$

Note that a special case of this yields the propositional uniqueness principle for products: $z = (\text{pr}_1(z), \text{pr}_2(z))$.

It can be helpful to view $\text{pair}^=$ as a *constructor* or *introduction rule* for $x = y$, analogous to the “pairing” constructor of $A \times B$ itself, which introduces the pair (a, b) given $a : A$ and $b : B$. From this perspective, the two components of (2.6.1):

$$\begin{aligned} \text{ap}_{\text{pr}_1} &: (x = y) \rightarrow (\text{pr}_1(x) = \text{pr}_1(y)) \\ \text{ap}_{\text{pr}_2} &: (x = y) \rightarrow (\text{pr}_2(x) = \text{pr}_2(y)) \end{aligned}$$

are *elimination* rules. Similarly, the two homotopies which witness (2.6.3) as quasi-inverse to (2.6.1) consist, respectively, of *propositional computation rules*:

$$\begin{aligned} \text{ap}_{\text{pr}_1}(\text{pair}^=(p, q)) &= p \\ \text{ap}_{\text{pr}_2}(\text{pair}^=(p, q)) &= q \end{aligned}$$

for $p : \text{pr}_1 x = \text{pr}_1 y$ and $q : \text{pr}_2 x = \text{pr}_2 y$, and a *propositional uniqueness principle*:

$$r = \text{pair}^=(\text{ap}_{\text{pr}_1}(r), \text{ap}_{\text{pr}_2}(r)) \quad \text{for } r : x =_{A \times B} y.$$

We can also characterize the reflexivity, inverses, and composition of paths in $A \times B$ componentwise:

$$\begin{aligned} \text{refl}_{(z:A \times B)} &= \text{pair}^=(\text{refl}_{\text{pr}_1 z}, \text{refl}_{\text{pr}_2 z}) \\ p^{-1} &= \text{pair}^=(\text{ap}_{\text{pr}_1}(p)^{-1}, \text{ap}_{\text{pr}_2}(p)^{-1}) \\ p \bullet q &= \text{pair}^=(\text{ap}_{\text{pr}_1}(p) \bullet \text{ap}_{\text{pr}_1}(q), \text{ap}_{\text{pr}_2}(p) \bullet \text{ap}_{\text{pr}_2}(q)). \end{aligned}$$

Or, written differently:

$$\begin{aligned} \text{ap}_{\text{pr}_i}(\text{refl}_{(z:A \times B)}) &= \text{refl}_{\text{pr}_i z} \quad (i = 1, 2) \\ \text{pair}^=(p^{-1}, q^{-1}) &= \text{pair}^=(p, q)^{-1} \\ \text{pair}^=(p \bullet q, p' \bullet q') &= \text{pair}^=(p, p') \bullet \text{pair}^=(q, q'). \end{aligned}$$

All of these equations can be derived by using path induction on the given paths and then returning reflexivity. The same is true for the rest of the higher groupoid structure considered in §2.1, although it begins to get tedious to insert enough other coherence paths to yield an equation that will typecheck. For instance, if we denote the inverse of the path in Lemma 2.1.4(iv) by

$\text{assoc}(p, q, r)$ and the last path displayed above by $\text{pair}^*(p, q, p', q')$, then for any $u, v, z, w : A \times B$ and p, q, r, p', q', r' of appropriate types we have

$$\begin{aligned} & \text{pair}^*(p \cdot q, r, p' \cdot q', r') \\ & \quad \bullet (\text{pair}^*(p, q, p', q') \cdot_r \text{pair}^=(r, r')) \\ & \quad \bullet \text{assoc}(\text{pair}^=(p, p'), \text{pair}^=(q, q'), \text{pair}^=(r, r')) \\ & = \text{ap}_{\text{pair}^=}\left(\text{pair}^=(\text{assoc}(p, q, r), \text{assoc}(p', q', r'))\right) \\ & \quad \bullet \text{pair}^*(p, q \cdot r, p', q' \cdot r') \\ & \quad \bullet (\text{pair}^=(p, p') \cdot_l \text{pair}^*(q, r, q', r')). \end{aligned}$$

Fortunately, we will never have to use any such higher-dimensional coherences.

We now consider transport in a pointwise product of type families. Given type families $A, B : Z \rightarrow \mathcal{U}$, we abusively write $A \times B : Z \rightarrow \mathcal{U}$ for the type family defined by $(A \times B)(z) := A(z) \times B(z)$. Now given $p : z =_Z w$ and $x : A(z) \times B(z)$, we can transport x along p to obtain an element of $A(w) \times B(w)$.

Theorem 2.6.4. *In the above situation, we have*

$$\text{transport}^{A \times B}(p, x) =_{A(w) \times B(w)} (\text{transport}^A(p, \text{pr}_1 x), \text{transport}^B(p, \text{pr}_2 x)).$$

Proof. By path induction, we may assume p is reflexivity, in which case we have

$$\begin{aligned} \text{transport}^{A \times B}(p, x) &\equiv x \\ \text{transport}^A(p, \text{pr}_1 x) &\equiv \text{pr}_1 x \\ \text{transport}^B(p, \text{pr}_2 x) &\equiv \text{pr}_2 x. \end{aligned}$$

Thus, it remains to show $x = (\text{pr}_1 x, \text{pr}_2 x)$. But this is the propositional uniqueness principle for product types, which, as we remarked above, follows from Theorem 2.6.2. \square

Finally, we consider the functoriality of ap under cartesian products. Suppose given types A, B, A', B' and functions $g : A \rightarrow A'$ and $h : B \rightarrow B'$; then we can define a function $f : A \times B \rightarrow A' \times B'$ by $f(x) := (g(\text{pr}_1 x), h(\text{pr}_2 x))$.

Theorem 2.6.5. *In the above situation, given $x, y : A \times B$ and $p : \text{pr}_1 x = \text{pr}_1 y$ and $q : \text{pr}_2 x = \text{pr}_2 y$, we have*

$$f(\text{pair}^=(p, q)) =_{(f(x)=f(y))} \text{pair}^=(g(p), h(q)).$$

Proof. Note first that the above equation is well-typed. On the one hand, since $\text{pair}^=(p, q) : x = y$ we have $f(\text{pair}^=(p, q)) : f(x) = f(y)$. On the other hand, since $\text{pr}_1(f(x)) \equiv g(\text{pr}_1 x)$ and $\text{pr}_2(f(x)) \equiv h(\text{pr}_2 x)$, we also have $\text{pair}^=(g(p), h(q)) : f(x) = f(y)$.

Now, by induction, we may assume $x \equiv (a, b)$ and $y \equiv (a', b')$, in which case we have $p : a = a'$ and $q : b = b'$. Thus, by path induction, we may assume p and q are reflexivity, in which case the desired equation holds judgmentally. \square

2.7 Σ -types

Let A be a type and $P : A \rightarrow \mathcal{U}$ a type family. Recall that the Σ -type, or dependent pair type, $\Sigma_{(x:A)} P(x)$ is a generalization of the cartesian product type. Thus, we expect its higher groupoid

structure to also be a generalization of the previous section. In particular, its paths should be pairs of paths, but it takes a little thought to give the correct types of these paths.

Suppose that we have a path $p : w = w'$ in $\sum_{(x:A)} P(x)$. Then we get $\text{pr}_1(p) : \text{pr}_1(w) = \text{pr}_1(w')$. However, we cannot directly ask whether $\text{pr}_2(w)$ is identical to $\text{pr}_2(w')$ since they don't have to be in the same type. But we can transport $\text{pr}_2(w)$ along the path $\text{pr}_1(p)$, and this does give us an element of the same type as $\text{pr}_2(w')$. By path induction, we do in fact obtain a path $\text{pr}_1(p)_*(\text{pr}_2(w)) = \text{pr}_2(w')$.

Recall from the discussion preceding Lemma 2.3.4 that $\text{pr}_1(p)_*(\text{pr}_2(w)) = \text{pr}_2(w')$ can be regarded as the type of paths from $\text{pr}_2(w)$ to $\text{pr}_2(w')$ which lie over the path $\text{pr}_1(p)$ in A . Thus, we are saying that a path $w = w'$ in the total space determines (and is determined by) a path $p : \text{pr}_1(w) = \text{pr}_1(w')$ in A together with a path from $\text{pr}_2(w)$ to $\text{pr}_2(w')$ lying over p , which seems sensible.

Remark 2.7.1. Note that if we have $x : A$ and $u, v : P(x)$ such that $(x, u) = (x, v)$, it does not follow that $u = v$. All we can conclude is that there exists $p : x = x$ such that $p_*(u) = v$. This is a well-known source of confusion for newcomers to type theory, but it makes sense from a topological viewpoint: the existence of a path $(x, u) = (x, v)$ in the total space of a fibration between two points that happen to lie in the same fiber does not imply the existence of a path $u = v$ lying entirely *within* that fiber.

The next theorem states that we can also reverse this process. Since it is a direct generalization of Theorem 2.6.2, we will be more concise.

Theorem 2.7.2. Suppose that $P : A \rightarrow \mathcal{U}$ is a type family over a type A and let $w, w' : \sum_{(x:A)} P(x)$. Then there is an equivalence

$$(w = w') \simeq \sum_{(p:\text{pr}_1(w) = \text{pr}_1(w'))} p_*(\text{pr}_2(w)) = \text{pr}_2(w').$$

Proof. We define a function

$$f : \prod_{w, w' : \sum_{(x:A)} P(x)} (w = w') \rightarrow \sum_{(p:\text{pr}_1(w) = \text{pr}_1(w'))} p_*(\text{pr}_2(w)) = \text{pr}_2(w')$$

by path induction, with

$$f(w, w, \text{refl}_w) := (\text{refl}_{\text{pr}_1(w)}, \text{refl}_{\text{pr}_2(w)}).$$

We want to show that f is an equivalence.

In the reverse direction, we define

$$g : \prod_{w, w' : \sum_{(x:A)} P(x)} \left(\sum_{p:\text{pr}_1(w) = \text{pr}_1(w')} p_*(\text{pr}_2(w)) = \text{pr}_2(w') \right) \rightarrow (w = w')$$

by first inducting on w and w' , which splits them into (w_1, w_2) and (w'_1, w'_2) respectively, so it suffices to show

$$\left(\sum_{p:w_1=w'_1} p_*(w_2) = w'_2 \right) \rightarrow ((w_1, w_2) = (w'_1, w'_2)).$$

Next, given a pair $\sum_{(p:w_1=w'_1)} p_*(w_2) = w'_2$, we can use Σ -induction to get $p : w_1 = w'_1$ and $q : p_*(w_2) = w'_2$. Inducting on p , we have $q : (\text{refl}_{w_1})_*(w_2) = w'_2$, and it suffices to show $(w_1, w_2) = (w_1, w'_2)$. But $(\text{refl}_{w_1})_*(w_2) \equiv w_2$, so inducting on q reduces the goal to $(w_1, w_2) = (w_1, w_2)$, which we can prove with $\text{refl}_{(w_1, w_2)}$.

Next we show that $f(g(r)) = r$ for all w, w' and r , where r has type

$$\sum_{(p:\text{pr}_1(w)=\text{pr}_1(w'))} (p_*(\text{pr}_2(w)) = \text{pr}_2(w')).$$

First, we break apart the pairs w, w' , and r by pair induction, as in the definition of g , and then use two path inductions to reduce both components of r to refl . Then it suffices to show that $f(g(\text{refl}_{w_1}, \text{refl}_{w_2})) = (\text{refl}_{w_1}, \text{refl}_{w_2})$, which is true by definition.

Similarly, to show that $g(f(p)) = p$ for all w, w' , and $p : w = w'$, we can do path induction on p , and then pair induction to split w , at which point it suffices to show that $g(f(\text{refl}_{(w_1, w_2)})) = \text{refl}_{(w_1, w_2)}$, which is true by definition.

Thus, f has a quasi-inverse, and is therefore an equivalence. \square

As we did in the case of cartesian products, we can deduce a propositional uniqueness principle as a special case.

Corollary 2.7.3. *For $z : \sum_{(x:A)} P(x)$, we have $z = (\text{pr}_1(z), \text{pr}_2(z))$.*

Proof. We have $\text{refl}_{\text{pr}_1(z)} : \text{pr}_1(z) = \text{pr}_1(\text{pr}_1(z), \text{pr}_2(z))$, so by Theorem 2.7.2 it will suffice to exhibit a path $(\text{refl}_{\text{pr}_1(z)})_* (\text{pr}_2(z)) = \text{pr}_2(\text{pr}_1(z), \text{pr}_2(z))$. But both sides are judgmentally equal to $\text{pr}_2(z)$. \square

Like with binary cartesian products, we can think of the backward direction of Theorem 2.7.2 as an introduction form ($\text{pair}^=$), the forward direction as elimination forms (ap_{pr_1} and ap_{pr_2}), and the equivalence as giving a propositional computation rule and uniqueness principle for these.

Note that the lifted path $\text{lift}(u, p)$ of $p : x = y$ at $u : P(x)$ defined in Lemma 2.3.2 may be identified with the special case of the introduction form

$$\text{pair}^= (p, \text{refl}_{p_*(u)}) : (x, u) = (y, p_*(u)).$$

This appears in the statement of action of transport on Σ -types, which is also a generalization of the action for binary cartesian products:

Theorem 2.7.4. *Suppose we have type families*

$$P : A \rightarrow \mathcal{U} \quad \text{and} \quad Q : \left(\sum_{x:A} P(x) \right) \rightarrow \mathcal{U}.$$

Then we can construct the type family over A defined by

$$x \mapsto \sum_{u:P(x)} Q(x, u).$$

For any path $p : x = y$ and any $(u, z) : \sum_{(u:P(x))} Q(x, u)$ we have

$$p_*(u, z) = (p_*(u), \text{pair}^= (p, \text{refl}_{p_*(u)})_*(z)).$$

Proof. Immediate by path induction. \square

We leave it to the reader to state and prove a generalization of Theorem 2.6.5 (see Exercise 2.7), and to characterize the reflexivity, inverses, and composition of Σ -types componentwise.

2.8 The unit type

Trivial cases are sometimes important, so we mention briefly the case of the unit type $\mathbf{1}$.

Theorem 2.8.1. *For any $x, y : \mathbf{1}$, we have $(x = y) \simeq \mathbf{1}$.*

It may be tempting to begin this proof by $\mathbf{1}$ -induction on x and y , reducing the problem to $(\star = \star) \simeq \mathbf{1}$. However, at this point we would be stuck, since we would be unable to perform a path induction on $p : \star = \star$. Thus, we instead work with a general x and y as much as possible, reducing them to \star by induction only at the last moment.

Proof. A function $(x = y) \rightarrow \mathbf{1}$ is easy to define by sending everything to \star . Conversely, for any $x, y : \mathbf{1}$ we may assume by induction that $x \equiv \star \equiv y$. In this case we have $\text{refl}_\star : x = y$, yielding a constant function $\mathbf{1} \rightarrow (x = y)$.

To show that these are inverses, consider first an element $u : \mathbf{1}$. We may assume that $u \equiv \star$, but this is also the result of the composite $\mathbf{1} \rightarrow (x = y) \rightarrow \mathbf{1}$.

On the other hand, suppose given $p : x = y$. By path induction, we may assume $x \equiv y$ and p is refl_x . We may then assume that x is \star , in which case the composite $(x = y) \rightarrow \mathbf{1} \rightarrow (x = y)$ takes p to refl_x , i.e. to p . \square

In particular, any two elements of $\mathbf{1}$ are equal. We leave it to the reader to formulate this equivalence in terms of introduction, elimination, computation, and uniqueness rules. The transport lemma for $\mathbf{1}$ is simply the transport lemma for constant type families (Lemma 2.3.5).

2.9 Π -types and the function extensionality axiom

Given a type A and a type family $B : A \rightarrow \mathcal{U}$, consider the dependent function type $\prod_{(x:A)} B(x)$. We expect the type $f = g$ of paths from f to g in $\prod_{(x:A)} B(x)$ to be equivalent to the type of pointwise paths:

$$(f = g) \simeq \left(\prod_{x:A} (f(x) =_{B(x)} g(x)) \right). \quad (2.9.1)$$

From a traditional perspective, this would say that two functions which are equal at each point are equal as functions. From a topological perspective, it would say that a path in a function space is the same as a continuous homotopy. And from a categorical perspective, it would say that an isomorphism in a functor category is a natural family of isomorphisms.

Unlike the case in the previous sections, however, the basic type theory presented in Chapter 1 is insufficient to prove (2.9.1). All we can say is that there is a certain function

$$\text{happly} : (f = g) \rightarrow \prod_{x:A} (f(x) =_{B(x)} g(x)) \quad (2.9.2)$$

which is easily defined by path induction. For the moment, therefore, we will assume:

Axiom 2.9.3 (Function extensionality). *For any A, B, f , and g , the function (2.9.2) is an equivalence.*

We will see in later chapters that this axiom follows both from univalence (see §§2.10 and 4.9) and from an interval type (see §6.3 and Exercise 6.10).

In particular, Axiom 2.9.3 implies that (2.9.2) has a quasi-inverse

$$\text{funext} : \left(\prod_{x:A} (f(x) = g(x)) \right) \rightarrow (f = g).$$

This function is also referred to as “function extensionality”. As we did with $\text{pair}^=$ in §2.6, we can regard funext as an *introduction rule* for the type $f = g$. From this point of view, happly is the *elimination rule*, while the homotopies witnessing funext as quasi-inverse to happly become a propositional computation rule

$$\text{happly}(\text{funext}(h), x) = h(x) \quad \text{for } h : \prod_{x:A} (f(x) = g(x))$$

and a propositional uniqueness principle:

$$p = \text{funext}(x \mapsto \text{happly}(p, x)) \quad \text{for } p : f = g.$$

We can also compute the identity, inverses, and composition in Π -types; they are simply given by pointwise operations::

$$\begin{aligned} \text{refl}_f &= \text{funext}(x \mapsto \text{refl}_{f(x)}) \\ \alpha^{-1} &= \text{funext}(x \mapsto \text{happly}(\alpha, x)^{-1}) \\ \alpha \bullet \beta &= \text{funext}(x \mapsto \text{happly}(\alpha, x) \bullet \text{happly}(\beta, x)). \end{aligned}$$

The first of these equalities follows from the definition of happly , while the second and third are easy path inductions.

Since the non-dependent function type $A \rightarrow B$ is a special case of the dependent function type $\prod_{(x:A)} B(x)$ when B is independent of x , everything we have said above applies in non-dependent cases as well. The rules for transport, however, are somewhat simpler in the non-dependent case. Given a type X , a path $p : x_1 =_X x_2$, type families $A, B : X \rightarrow \mathcal{U}$, and a function $f : A(x_1) \rightarrow B(x_1)$, we have

$$\text{transport}^{A \rightarrow B}(p, f) = \left(x \mapsto \text{transport}^B(p, f(\text{transport}^A(p^{-1}, x))) \right) \quad (2.9.4)$$

where $A \rightarrow B$ denotes abusively the type family $X \rightarrow \mathcal{U}$ defined by

$$(A \rightarrow B)(x) := (A(x) \rightarrow B(x)).$$

In other words, when we transport a function $f : A(x_1) \rightarrow B(x_1)$ along a path $p : x_1 = x_2$, we obtain the function $A(x_2) \rightarrow B(x_2)$ which transports its argument backwards along p (in the type family A), applies f , and then transports the result forwards along p (in the type family B). This can be proven easily by path induction.

Transporting dependent functions is similar, but more complicated. Suppose given X and p as before, type families $A : X \rightarrow \mathcal{U}$ and $B : \prod_{(x:X)} (A(x) \rightarrow \mathcal{U})$, and also a dependent function $f : \prod_{(a:A(x_1))} B(x_1, a)$. Then for $a : A(x_2)$, we have

$$\text{transport}^{\prod_A B}(p, f)(a) = \text{transport}^{\widehat{B}} \left((\text{pair}^=(p^{-1}, \text{refl}_{p^{-1}*a}))^{-1}, f(\text{transport}^A(p^{-1}, a)) \right)$$

where $\prod_A B$ and \widehat{B} denote respectively the type families

$$\begin{aligned} \prod_A B &:= (x \mapsto \prod_{(a:A(x))} B(x, a)) : X \rightarrow \mathcal{U} \\ \widehat{B} &:= (w \mapsto B(\text{pr}_1 w, \text{pr}_2 w)) : (\sum_{(x:X)} A(x)) \rightarrow \mathcal{U}. \end{aligned} \quad (2.9.5)$$

If these formulas look a bit intimidating, don’t worry about the details. The basic idea is just the same as for the non-dependent function type: we transport the argument backwards, apply the function, and then transport the result forwards again.

Now recall that for a general type family $P : X \rightarrow \mathcal{U}$, in §2.2 we defined the type of *dependent paths* over $p : x =_X y$ from $u : P(x)$ to $v : P(y)$ to be $p_*(u) =_{P(y)} v$. When P is a family of function types, there is an equivalent way to represent this which is often more convenient.

Lemma 2.9.6. *Given type families $A, B : X \rightarrow \mathcal{U}$ and $p : x =_X y$, and also $f : A(x) \rightarrow B(x)$ and $g : A(y) \rightarrow B(y)$, we have an equivalence*

$$(p_*(f) = g) \simeq \prod_{a:A(x)} (p_*(f(a)) = g(p_*(a))).$$

Moreover, if $q : p_*(f) = g$ corresponds under this equivalence to \hat{q} , then for $a : A(x)$, the path

$$\text{happly}(q, p_*(a)) : (p_*(f))(p_*(a)) = g(p_*(a))$$

is equal to the concatenated path $i \cdot j \cdot k$, where

- $i : (p_*(f))(p_*(a)) = p_*(f(p^{-1}_*(p_*(a))))$ comes from (2.9.4),
- $j : p_*(f(p^{-1}_*(p_*(a)))) = p_*(f(a))$ comes from Lemmas 2.1.4 and 2.3.9, and
- $k : p_*(f(a)) = g(p_*(a))$ is $\hat{q}(a)$.

Proof. By path induction, we may assume p is reflexivity, in which case the desired equivalence reduces to function extensionality. The second statement then follows by the computation rule for function extensionality. \square

In general, it happens quite frequently that we want to consider a concatenation of paths each of which arises from some previously proven lemmas or hypothesized objects, and it can be rather tedious to describe this by giving a name to each path in the concatenation as we did in the second statement above. Thus, we adopt a convention of writing such concatenations in the familiar mathematical style of “chains of equalities with reasons”, and allow ourselves to omit reasons that the reader can easily fill in. For instance, the path $i \cdot j \cdot k$ from Lemma 2.9.6 would be written like this:

$$\begin{aligned} (p_*(f))(p_*(a)) &= p_*\left(f(p^{-1}_*(p_*(a)))\right) && \text{(by (2.9.4))} \\ &= p_*(f(a)) \\ &= g(p_*(a)). && \text{(by } \hat{q} \text{)} \end{aligned}$$

In ordinary mathematics, such a chain of equalities would be merely proving that two things are equal. We are enhancing this by using it to describe a *particular* path between them.

As usual, there is a version of Lemma 2.9.6 for dependent functions that is similar, but more complicated.

Lemma 2.9.7. *Given type families $A : X \rightarrow \mathcal{U}$ and $B : \prod_{(x:X)} A(x) \rightarrow \mathcal{U}$ and $p : x =_X y$, and also $f : \prod_{(a:A(x))} B(x, a)$ and $g : \prod_{(a:A(y))} B(y, a)$, we have an equivalence*

$$(p_*(f) = g) \simeq \left(\prod_{a:A(x)} \text{transport}^{\hat{B}}(\text{pair}^=(p, \text{refl}_{p_*(a)}), f(a)) = g(p_*(a)) \right)$$

with \hat{B} as in (2.9.5).

We leave it to the reader to prove this and to formulate a suitable computation rule.

2.10 Universes and the univalence axiom

Given two types A and B , we may consider them as elements of some universe type \mathcal{U} , and thereby form the identity type $A =_{\mathcal{U}} B$. As mentioned in the introduction, *univalence* is the identification of $A =_{\mathcal{U}} B$ with the type $(A \simeq B)$ of equivalences from A to B , which we described in §2.4. We perform this identification by way of the following canonical function.

Lemma 2.10.1. *For types $A, B : \mathcal{U}$, there is a certain function,*

$$\text{idtoeqv} : (A =_{\mathcal{U}} B) \rightarrow (A \simeq B), \quad (2.10.2)$$

defined in the proof.

Proof. We could construct this directly by induction on equality, but the following description is more convenient. Note that the identity function $\text{id}_{\mathcal{U}} : \mathcal{U} \rightarrow \mathcal{U}$ may be regarded as a type family indexed by the universe \mathcal{U} ; it assigns to each type $X : \mathcal{U}$ the type X itself. (When regarded as a fibration, its total space is the type $\sum_{(A:\mathcal{U})} A$ of “pointed types”; see also §4.8.) Thus, given a path $p : A =_{\mathcal{U}} B$, we have a transport function $p_* : A \rightarrow B$. We claim that p_* is an equivalence. But by induction, it suffices to assume that p is refl_A , in which case $p_* \equiv \text{id}_A$, which is an equivalence by Example 2.4.7. Thus, we can define $\text{idtoeqv}(p)$ to be p_* (together with the above proof that it is an equivalence). \square

We would like to say that idtoeqv is an equivalence. However, as with happly for function types, the type theory described in Chapter 1 is insufficient to guarantee this. Thus, as we did for function extensionality, we formulate this property as an axiom: Voevodsky’s *univalence axiom*.

Axiom 2.10.3 (Univalence). *For any $A, B : \mathcal{U}$, the function (2.10.2) is an equivalence.*

In particular, therefore, we have

$$(A =_{\mathcal{U}} B) \simeq (A \simeq B).$$

Technically, the univalence axiom is a statement about a particular universe type \mathcal{U} . If a universe \mathcal{U} satisfies this axiom, we say that it is **univalent**. Except when otherwise noted (e.g. in §4.9) we will assume that *all* universes are univalent.

Remark 2.10.4. It is important for the univalence axiom that we defined $A \simeq B$ using a “good” version of isequiv as described in §2.4, rather than (say) as $\sum_{(f:A \rightarrow B)} \text{qinv}(f)$. See Exercise 4.6.

In particular, univalence means that *equivalent types may be identified*. As we did in previous sections, it is useful to break this equivalence into:

- An introduction rule for $(A =_{\mathcal{U}} B)$, denoted ua for “univalence axiom”:

$$\text{ua} : (A \simeq B) \rightarrow (A =_{\mathcal{U}} B).$$

- The elimination rule, which is idtoeqv ,

$$\text{idtoeqv} \equiv \text{transport}^{X \mapsto X} : (A =_{\mathcal{U}} B) \rightarrow (A \simeq B).$$

- The propositional computation rule,

$$\text{transport}^{X \mapsto X}(\text{ua}(f), x) = f(x).$$

- The propositional uniqueness principle: for any $p : A = B$,

$$p = \text{ua}(\text{transport}^{X \hookrightarrow X}(p)).$$

We can also identify the reflexivity, concatenation, and inverses of equalities in the universe with the corresponding operations on equivalences:

$$\begin{aligned}\text{refl}_A &= \text{ua}(\text{id}_A) \\ \text{ua}(f) \bullet \text{ua}(g) &= \text{ua}(g \circ f) \\ \text{ua}(f)^{-1} &= \text{ua}(f^{-1}).\end{aligned}$$

The first of these follows because $\text{id}_A = \text{idtoeqv}(\text{refl}_A)$ by definition of idtoeqv , and ua is the inverse of idtoeqv . For the second, if we define $p := \text{ua}(f)$ and $q := \text{ua}(g)$, then we have

$$\text{ua}(g \circ f) = \text{ua}(\text{idtoeqv}(q) \circ \text{idtoeqv}(p)) = \text{ua}(\text{idtoeqv}(p \bullet q)) = p \bullet q$$

using Lemma 2.3.9 and the definition of idtoeqv . The third is similar.

The following observation, which is a special case of Lemma 2.3.10, is often useful when applying the univalence axiom.

Lemma 2.10.5. *For any type family $B : A \rightarrow \mathcal{U}$ and $x, y : A$ with a path $p : x = y$ and $u : B(x)$, we have*

$$\begin{aligned}\text{transport}^B(p, u) &= \text{transport}^{X \hookrightarrow X}(\text{ap}_B(p), u) \\ &= \text{idtoeqv}(\text{ap}_B(p))(u).\end{aligned}$$

2.11 Identity type

Just as the type $a =_A a'$ is characterized up to isomorphism, with a separate “definition” for each A , there is no simple characterization of the type $p =_{a=_A a'} q$ of paths between paths $p, q : a =_A a'$. However, our other general classes of theorems do extend to identity types, such as the fact that they respect equivalence.

Theorem 2.11.1. *If $f : A \rightarrow B$ is an equivalence, then for all $a, a' : A$, so is*

$$\text{ap}_f : (a =_A a') \rightarrow (f(a) =_B f(a')).$$

Proof. Let f^{-1} be a quasi-inverse of f , with homotopies

$$\alpha : \prod_{b:B} (f(f^{-1}(b)) = b) \quad \text{and} \quad \beta : \prod_{a:A} (f^{-1}(f(a)) = a).$$

The quasi-inverse of ap_f is, essentially,

$$\text{ap}_{f^{-1}} : (f(a) = f(a')) \rightarrow (f^{-1}(f(a)) = f^{-1}(f(a'))).$$

However, in order to obtain an element of $a =_A a'$ from $\text{ap}_{f^{-1}}(q)$, we must concatenate with the paths β_a^{-1} and $\beta_{a'}$ on either side. To show that this gives a quasi-inverse of ap_f , on one hand we must show that for any $p : a =_A a'$ we have

$$\beta_a^{-1} \bullet \text{ap}_{f^{-1}}(\text{ap}_f(p)) \bullet \beta_{a'} = p.$$

This follows from the functoriality of ap and the naturality of homotopies, Lemmas 2.2.2 and 2.4.3. On the other hand, we must show that for any $q : f(a) =_B f(a')$ we have

$$\text{ap}_f(\beta_a^{-1} \cdot \text{ap}_{f^{-1}}(q) \cdot \beta_{a'}) = q.$$

The proof of this is a little more involved, but each step is again an application of Lemmas 2.2.2 and 2.4.3 (or simply canceling inverse paths):

$$\begin{aligned} \text{ap}_f(\beta_a^{-1} \cdot \text{ap}_{f^{-1}}(q) \cdot \beta_{a'}) &= \alpha_{f(a)}^{-1} \cdot \alpha_{f(a)} \cdot \text{ap}_f(\beta_a^{-1} \cdot \text{ap}_{f^{-1}}(q) \cdot \beta_{a'}) \cdot \alpha_{f(a')}^{-1} \cdot \alpha_{f(a')} \\ &= \alpha_{f(a)}^{-1} \cdot \text{ap}_f(\text{ap}_{f^{-1}}(\text{ap}_f(\beta_a^{-1} \cdot \text{ap}_{f^{-1}}(q) \cdot \beta_{a'}))) \cdot \alpha_{f(a')} \\ &= \alpha_{f(a)}^{-1} \cdot \text{ap}_f(\beta_a \cdot \beta_a^{-1} \cdot \text{ap}_{f^{-1}}(q) \cdot \beta_{a'} \cdot \beta_{a'}^{-1}) \cdot \alpha_{f(a')} \\ &= \alpha_{f(a)}^{-1} \cdot \text{ap}_f(\text{ap}_{f^{-1}}(q)) \cdot \alpha_{f(a')} \\ &= q. \end{aligned} \quad \square$$

Thus, if for some type A we have a full characterization of $a =_A a'$, the type $p =_{a=_A a'} q$ is determined as well. For example:

- Paths $p = q$, where $p, q : w =_{A \times B} w'$, are equivalent to pairs of paths

$$\text{ap}_{\text{pr}_1} p =_{\text{pr}_1 w =_A \text{pr}_1 w'} \text{ap}_{\text{pr}_1} q \quad \text{and} \quad \text{ap}_{\text{pr}_2} p =_{\text{pr}_2 w =_B \text{pr}_2 w'} \text{ap}_{\text{pr}_2} q.$$

- Paths $p = q$, where $p, q : f =_{\prod_{(x:A)} B(x)} g$, are equivalent to homotopies

$$\prod_{x:A} (\text{happly}(p)(x) =_{f(x)=g(x)} \text{happly}(q)(x)).$$

Next we consider transport in families of paths, i.e. transport in $C : A \rightarrow \mathcal{U}$ where each $C(x)$ is an identity type. The simplest case is when $C(x)$ is a type of paths in A itself, perhaps with one endpoint fixed.

Lemma 2.11.2. *For any A and $a : A$, with $p : x_1 = x_2$, we have*

$$\begin{aligned} \text{transport}^{x \mapsto (a=x)}(p, q) &= q \cdot p && \text{for } q : a = x_1, \\ \text{transport}^{x \mapsto (x=a)}(p, q) &= p^{-1} \cdot q && \text{for } q : x_1 = a, \\ \text{transport}^{x \mapsto (x=x)}(p, q) &= p^{-1} \cdot q \cdot p && \text{for } q : x_1 = x_1. \end{aligned}$$

Proof. Path induction on p , followed by the unit laws for composition. \square

In other words, transporting with $x \mapsto c = x$ is post-composition, and transporting with $x \mapsto x = c$ is contravariant pre-composition. These may be familiar as the functorial actions of the covariant and contravariant hom-functors $\text{hom}(c, -)$ and $\text{hom}(-, c)$ in category theory.

Similarly, we can prove the following more general form of Lemma 2.11.2, which is related to Lemma 2.3.10.

Theorem 2.11.3. *For $f, g : A \rightarrow B$, with $p : a =_A a'$ and $q : f(a) =_B g(a)$, we have*

$$\text{transport}^{x \mapsto f(x)=_B g(x)}(p, q) =_{f(a')=g(a')} (\text{ap}_f p)^{-1} \cdot q \cdot \text{ap}_g p.$$

Because $\text{ap}_{(x \mapsto x)}$ is the identity function and $\text{ap}_{(x \mapsto c)}$ (where c is a constant) is $p \mapsto \text{refl}_c$, Lemma 2.11.2 is a special case. A yet more general version is when B can be a family of types indexed on A :

Theorem 2.11.4. Let $B : A \rightarrow \mathcal{U}$ and $f, g : \prod_{(x:A)} B(x)$, with $p : a =_A a'$ and $q : f(a) =_{B(a)} g(a)$. Then we have

$$\text{transport}^{x \mapsto f(x) =_{B(x)} g(x)}(p, q) = (\text{apd}_f(p))^{-1} \cdot \text{ap}_{(\text{transport}^B p)}(q) \cdot \text{apd}_g(p).$$

Finally, as in §2.9, for families of identity types there is another equivalent characterization of dependent paths.

Theorem 2.11.5. For $p : a =_A a'$ with $q : a = a$ and $r : a' = a'$, we have

$$(\text{transport}^{x \mapsto (x=x)}(p, q) = r) \simeq (q \cdot p = p \cdot r).$$

Proof. Path induction on p , followed by the fact that composing with the unit equalities $q \cdot 1 = q$ and $r = 1 \cdot r$ is an equivalence. \square

There are more general equivalences involving the application of functions, akin to Theorems 2.11.3 and 2.11.4.

2.12 Coproducts

So far, most of the type formers we have considered have been what are called *negative*. Intuitively, this means that their elements are determined by their behavior under the elimination rules: a (dependent) pair is determined by its projections, and a (dependent) function is determined by its values. The identity types of negative types can almost always be characterized straightforwardly, along with all of their higher structure, as we have done in §§2.6–2.9. The universe is not exactly a negative type, but its identity types behave similarly: we have a straightforward characterization (univalence) and a description of the higher structure. Identity types themselves, of course, are a special case.

We now consider our first example of a *positive* type former. Again informally, a positive type is one which is “presented” by certain constructors, with the universal property of a presentation being expressed by its elimination rule. (Categorically speaking, a positive type has a “mapping out” universal property, while a negative type has a “mapping in” universal property.) Because computing with presentations is, in general, an uncomputable problem, for positive types we cannot always expect a straightforward characterization of the identity type. However, in many particular cases, a characterization or partial characterization does exist, and can be obtained by the general method that we introduce with this example.

(Technically, our chosen presentation of cartesian products and Σ -types is also positive. However, because these types also admit a negative presentation which differs only slightly, their identity types have a direct characterization that does not require the method to be described here.)

Consider the coproduct type $A + B$, which is “presented” by the injections $\text{inl} : A \rightarrow A + B$ and $\text{inr} : B \rightarrow A + B$. Intuitively, we expect that $A + B$ contains exact copies of A and B disjointly, so that we should have

$$(\text{inl}(a_1) = \text{inl}(a_2)) \simeq (a_1 = a_2) \tag{2.12.1}$$

$$(\text{inr}(b_1) = \text{inr}(b_2)) \simeq (b_1 = b_2) \tag{2.12.2}$$

$$(\text{inl}(a) = \text{inr}(b)) \simeq \mathbf{0}. \tag{2.12.3}$$

We prove this as follows. Fix an element $a_0 : A$; we will characterize the type family

$$(x \mapsto (\text{inl}(a_0) = x)) : A + B \rightarrow \mathcal{U}. \quad (2.12.4)$$

A similar argument would characterize the analogous family $x \mapsto (x = \text{inr}(b_0))$ for any $b_0 : B$. Together, these characterizations imply (2.12.1)–(2.12.3).

In order to characterize (2.12.4), we will define a type family $\text{code} : A + B \rightarrow \mathcal{U}$ and show that $\prod_{(x:A+B)}((\text{inl}(a_0) = x) \simeq \text{code}(x))$. Since we want to conclude (2.12.1) from this, we should have $\text{code}(\text{inl}(a)) = (a_0 = a)$, and since we also want to conclude (2.12.3), we should have $\text{code}(\text{inr}(b)) = \mathbf{0}$. The essential insight is that we can use the recursion principle of $A + B$ to define $\text{code} : A + B \rightarrow \mathcal{U}$ by these two equations:

$$\begin{aligned} \text{code}(\text{inl}(a)) &:= (a_0 = a), \\ \text{code}(\text{inr}(b)) &:= \mathbf{0}. \end{aligned}$$

This is a very simple example of a proof technique that is used quite a bit when doing homotopy theory in homotopy type theory; see e.g. [????](#). We can now show:

Theorem 2.12.5. *For all $x : A + B$ we have $(\text{inl}(a_0) = x) \simeq \text{code}(x)$.*

Proof. The key to the following proof is that we do it for all points x together, enabling us to use the elimination principle for the coproduct. We first define a function

$$\text{encode} : \prod_{(x:A+B)} \prod_{(p:\text{inl}(a_0)=x)} \text{code}(x)$$

by transporting reflexivity along p :

$$\text{encode}(x, p) := \text{transport}^{\text{code}}(p, \text{refl}_{a_0}).$$

Note that $\text{refl}_{a_0} : \text{code}(\text{inl}(a_0))$, since $\text{code}(\text{inl}(a_0)) \equiv (a_0 = a_0)$ by definition of code . Next, we define a function

$$\text{decode} : \prod_{(x:A+B)} \prod_{(c:\text{code}(x))} (\text{inl}(a_0) = x).$$

To define $\text{decode}(x, c)$, we may first use the elimination principle of $A + B$ to divide into cases based on whether x is of the form $\text{inl}(a)$ or the form $\text{inr}(b)$.

In the first case, where $x \equiv \text{inl}(a)$, then $\text{code}(x) \equiv (a_0 = a)$, so that c is an identification between a_0 and a . Thus, $\text{ap}_{\text{inl}}(c) : (\text{inl}(a_0) = \text{inl}(a))$ so we can define this to be $\text{decode}(\text{inl}(a), c)$.

In the second case, where $x \equiv \text{inr}(b)$, then $\text{code}(x) \equiv \mathbf{0}$, so that c inhabits the empty type. Thus, the elimination rule of $\mathbf{0}$ yields a value for $\text{decode}(\text{inr}(b), c)$.

This completes the definition of decode ; we now show that $\text{encode}(x, -)$ and $\text{decode}(x, -)$ are quasi-inverses for all x . On the one hand, suppose given $x : A + B$ and $p : \text{inl}(a_0) = x$; we want to show $\text{decode}(x, \text{encode}(x, p)) = p$. But now by (based) path induction, it suffices to consider $x \equiv \text{inl}(a_0)$ and $p \equiv \text{refl}_{\text{inl}(a_0)}$:

$$\begin{aligned} \text{decode}(x, \text{encode}(x, p)) &\equiv \text{decode}(\text{inl}(a_0), \text{encode}(\text{inl}(a_0), \text{refl}_{\text{inl}(a_0)})) \\ &\equiv \text{decode}(\text{inl}(a_0), \text{transport}^{\text{code}}(\text{refl}_{\text{inl}(a_0)}, \text{refl}_{a_0})) \\ &\equiv \text{decode}(\text{inl}(a_0), \text{refl}_{a_0}) \\ &\equiv \text{inl}(\text{refl}_{a_0}) \\ &\equiv \text{refl}_{\text{inl}(a_0)} \\ &\equiv p. \end{aligned}$$

On the other hand, let $x : A + B$ and $c : \text{code}(x)$; we want to show $\text{encode}(x, \text{decode}(x, c)) = c$. We may again divide into cases based on x . If $x \equiv \text{inl}(a)$, then $c : a_0 = a$ and $\text{decode}(x, c) \equiv \text{ap}_{\text{inl}}(c)$, so that

$$\begin{aligned} \text{encode}(x, \text{decode}(x, c)) &\equiv \text{transport}^{\text{code}}(\text{ap}_{\text{inl}}(c), \text{refl}_{a_0}) \\ &= \text{transport}^{a \mapsto (a_0 = a)}(c, \text{refl}_{a_0}) && (\text{by Lemma 2.3.10}) \\ &= \text{refl}_{a_0} \bullet c && (\text{by Lemma 2.11.2}) \\ &= c. \end{aligned}$$

Finally, if $x \equiv \text{inr}(b)$, then $c : \mathbf{0}$, so we may conclude anything we wish. \square

Of course, there is a corresponding theorem if we fix $b_0 : B$ instead of $a_0 : A$.

In particular, Theorem 2.12.5 implies that for any $a : A$ and $b : B$ there are functions

$$\text{encode}(\text{inl}(a), -) : (\text{inl}(a_0) = \text{inl}(a)) \rightarrow (a_0 = a)$$

and

$$\text{encode}(\text{inr}(b), -) : (\text{inr}(b_0) = \text{inr}(b)) \rightarrow \mathbf{0}.$$

The second of these states “ $\text{inl}(a_0)$ is not equal to $\text{inr}(b)$ ”, i.e. the images of inl and inr are disjoint. The traditional reading of the first one, where identity types are viewed as propositions, is just injectivity of inl . The full homotopical statement of Theorem 2.12.5 gives more information: the types $\text{inl}(a_0) = \text{inl}(a)$ and $a_0 = a$ are actually equivalent, as are $\text{inr}(b_0) = \text{inr}(b)$ and $b_0 = b$.

Remark 2.12.6. In particular, since the two-element type $\mathbf{2}$ is equivalent to $\mathbf{1} + \mathbf{1}$, we have $0_2 \neq 1_2$.

This proof illustrates a general method for describing path spaces, which we will use often. To characterize a path space, the first step is to define a comparison fibration “code” that provides a more explicit description of the paths. There are several different methods for proving that such a comparison fibration is equivalent to the paths (we show a few different proofs of the same result in ??). The one we have used here is called the **encode-decode method**: the key idea is to define decode generally for all instances of the fibration (i.e. as a function $\prod_{(x:A+B)} \text{code}(x) \rightarrow (\text{inl}(a_0) = x)$), so that path induction can be used to analyze $\text{decode}(x, \text{encode}(x, p))$.

As usual, we can also characterize the action of transport in coproduct types. Given a type X , a path $p : x_1 =_X x_2$, and type families $A, B : X \rightarrow \mathcal{U}$, we have

$$\begin{aligned} \text{transport}^{A+B}(p, \text{inl}(a)) &= \text{inl}(\text{transport}^A(p, a)), \\ \text{transport}^{A+B}(p, \text{inr}(b)) &= \text{inr}(\text{transport}^B(p, b)), \end{aligned}$$

where as usual, $A + B$ in the superscript denotes abusively the type family $x \mapsto A(x) + B(x)$. The proof is an easy path induction.

2.13 Natural numbers

We use the encode-decode method to characterize the path space of the natural numbers, which are also a positive type. In this case, rather than fixing one endpoint, we characterize the two-sided path space all at once. Thus, the codes for identities are a type family

$$\text{code} : \mathbb{N} \rightarrow \mathbb{N} \rightarrow \mathcal{U},$$

defined by double recursion over \mathbb{N} as follows:

$$\begin{aligned} \text{code}(0, 0) &:= \mathbf{1} \\ \text{code}(\text{succ}(m), 0) &:= \mathbf{0} \\ \text{code}(0, \text{succ}(n)) &:= \mathbf{0} \\ \text{code}(\text{succ}(m), \text{succ}(n)) &:= \text{code}(m, n). \end{aligned}$$

We also define by recursion a dependent function $r : \prod_{(n:\mathbb{N})} \text{code}(n, n)$, with

$$\begin{aligned} r(0) &:= \star \\ r(\text{succ}(n)) &:= r(n). \end{aligned}$$

Theorem 2.13.1. *For all $m, n : \mathbb{N}$ we have $(m = n) \simeq \text{code}(m, n)$.*

Proof. We define

$$\text{encode} : \prod_{m,n:\mathbb{N}} (m = n) \rightarrow \text{code}(m, n)$$

by transporting, $\text{encode}(m, n, p) := \text{transport}^{\text{code}(m, -)}(p, r(m))$. And we define

$$\text{decode} : \prod_{m,n:\mathbb{N}} \text{code}(m, n) \rightarrow (m = n)$$

by double induction on m, n . When m and n are both 0, we need a function $\mathbf{1} \rightarrow (0 = 0)$, which we define to send everything to refl_0 . When m is a successor and n is 0 or vice versa, the domain $\text{code}(m, n)$ is $\mathbf{0}$, so the eliminator for $\mathbf{0}$ suffices. And when both are successors, we can define $\text{decode}(\text{succ}(m), \text{succ}(n))$ to be the composite

$$\text{code}(\text{succ}(m), \text{succ}(n)) \equiv \text{code}(m, n) \xrightarrow{\text{decode}(m, n)} (m = n) \xrightarrow{\text{ap}_{\text{succ}}} (\text{succ}(m) = \text{succ}(n)).$$

Next we show that $\text{encode}(m, n)$ and $\text{decode}(m, n)$ are quasi-inverses for all m, n .

On one hand, if we start with $p : m = n$, then by induction on p it suffices to show

$$\text{decode}(n, n, \text{encode}(n, n, \text{refl}_n)) = \text{refl}_n.$$

But $\text{encode}(n, n, \text{refl}_n) \equiv r(n)$, so it suffices to show that $\text{decode}(n, n, r(n)) = \text{refl}_n$. We can prove this by induction on n . If $n \equiv 0$, then $\text{decode}(0, 0, r(0)) = \text{refl}_0$ by definition of decode . And in the case of a successor, by the inductive hypothesis we have $\text{decode}(n, n, r(n)) = \text{refl}_n$, so it suffices to observe that $\text{ap}_{\text{succ}}(\text{refl}_n) \equiv \text{refl}_{\text{succ}(n)}$.

On the other hand, if we start with $c : \text{code}(m, n)$, then we proceed by double induction on m and n . If both are 0, then $\text{decode}(0, 0, c) \equiv \text{refl}_0$, while $\text{encode}(0, 0, \text{refl}_0) \equiv r(0) \equiv \star$. Thus, it suffices to recall from §2.8 that every inhabitant of $\mathbf{1}$ is equal to \star . If m is 0 but n is a successor, or vice versa, then $c : \mathbf{0}$, so we are done. And in the case of two successors, we have

$$\begin{aligned} \text{encode}(\text{succ}(m), \text{succ}(n), \text{decode}(\text{succ}(m), \text{succ}(n), c)) &= \text{encode}(\text{succ}(m), \text{succ}(n), \text{ap}_{\text{succ}}(\text{decode}(m, n, c))) \\ &= \text{transport}^{\text{code}(\text{succ}(m), -)}(\text{ap}_{\text{succ}}(\text{decode}(m, n, c)), r(\text{succ}(m))) \\ &= \text{transport}^{\text{code}(\text{succ}(m), \text{succ}(-))}(\text{decode}(m, n, c), r(\text{succ}(m))) \\ &= \text{transport}^{\text{code}(m, -)}(\text{decode}(m, n, c), r(m)) \\ &= \text{encode}(m, n, \text{decode}(m, n, c)) \\ &= c \end{aligned}$$

using the inductive hypothesis. □

In particular, we have

$$\text{encode}(\text{succ}(m), 0) : (\text{succ}(m) = 0) \rightarrow \mathbf{0} \quad (2.13.2)$$

which shows that “0 is not the successor of any natural number”. We also have the composite

$$(\text{succ}(m) = \text{succ}(n)) \xrightarrow{\text{encode}} \text{code}(\text{succ}(m), \text{succ}(n)) \equiv \text{code}(m, n) \xrightarrow{\text{decode}} (m = n) \quad (2.13.3)$$

which shows that the function succ is injective.

We will study more general positive types in Chapters 5 and 6. In ??, we will see that the same technique used here to characterize the identity types of coproducts and \mathbb{N} can also be used to calculate homotopy groups of spheres.

2.14 Example: equality of structures

We now consider one example to illustrate the interaction between the groupoid structure on a type and the type formers. In the introduction we remarked that one of the advantages of univalence is that two isomorphic things are interchangeable, in the sense that every property or construction involving one also applies to the other. Common “abuses of notation” become formally true. Univalence itself says that equivalent types are equal, and therefore interchangeable, which includes e.g. the common practice of identifying isomorphic sets. Moreover, when we define other mathematical objects as sets, or even general types, equipped with structure or properties, we can derive the correct notion of equality for them from univalence. We will illustrate this point with a significant example in ??, where we define the basic notions of category theory in such a way that equality of categories is equivalence, equality of functors is natural isomorphism, etc. See in particular ?? . In this section, we describe a very simple example, coming from algebra.

For simplicity, we use *semigroups* as our example, where a semigroup is a type equipped with an associative “multiplication” operation. The same ideas apply to other algebraic structures, such as monoids, groups, and rings. Recall from §§1.6 and 1.11 that the definition of a kind of mathematical structure should be interpreted as defining the type of such structures as a certain iterated Σ -type. In the case of semigroups this yields the following.

Definition 2.14.1. Given a type A , the type $\text{SemigroupStr}(A)$ of **semigroup structures** with carrier A is defined by

$$\text{SemigroupStr}(A) := \sum_{(m:A \rightarrow A \rightarrow A)} \prod_{(x,y,z:A)} m(x, m(y, z)) = m(m(x, y), z).$$

A **semigroup** is a type together with such a structure:

$$\text{Semigroup} := \sum_{A:\mathcal{U}} \text{SemigroupStr}(A)$$

In the next two sections, we describe two ways in which univalence makes it easier to work with such semigroups.

2.14.1 Lifting equivalences

When working loosely, one might say that a bijection between sets A and B “obviously” induces an isomorphism between semigroup structures on A and semigroup structures on B . With univalence, this is indeed obvious, because given an equivalence between types A and B , we can automatically derive a semigroup structure on B from one on A , and moreover show that this derivation is an equivalence of semigroup structures. The reason is that SemigroupStr is a family of types, and therefore has an action on paths between types given by transport:

$$\text{transport}^{\text{SemigroupStr}}(\text{ua}(e)) : \text{SemigroupStr}(A) \rightarrow \text{SemigroupStr}(B).$$

Moreover, this map is an equivalence, because $\text{transport}^C(\alpha)$ is always an equivalence with inverse $\text{transport}^C(\alpha^{-1})$, see Lemmas 2.1.4 and 2.3.9.

While the univalence axiom ensures that this map exists, we need to use facts about transport proven in the preceding sections to calculate what it actually does. Let (m, a) be a semigroup structure on A , and we investigate the induced semigroup structure on B given by

$$\text{transport}^{\text{SemigroupStr}}(\text{ua}(e), (m, a)).$$

First, because $\text{SemigroupStr}(X)$ is defined to be a Σ -type, by Theorem 2.7.4,

$$\text{transport}^{\text{SemigroupStr}}(\text{ua}(e), (m, a)) = (m', a')$$

where m' is an induced multiplication operation on B

$$\begin{aligned} m' &: B \rightarrow B \rightarrow B \\ m'(b_1, b_2) &\coloneqq \text{transport}^{X \mapsto (X \rightarrow X \rightarrow X)}(\text{ua}(e), m)(b_1, b_2) \end{aligned}$$

and a' an induced proof that m' is associative. We have, again by Theorem 2.7.4,

$$\begin{aligned} a' &: \text{Assoc}(B, m') \\ a' &\coloneqq \text{transport}^{(X, m) \mapsto \text{Assoc}(X, m)}((\text{pair}^=(\text{ua}(e), \text{refl}_{m'})), a), \end{aligned} \tag{2.14.2}$$

where $\text{Assoc}(X, m)$ is the type $\prod_{(x, y, z: X)} m(x, m(y, z)) = m(m(x, y), z)$. By function extensionality, it suffices to investigate the behavior of m' when applied to arguments $b_1, b_2 : B$. By applying (2.9.4) twice, we have that $m'(b_1, b_2)$ is equal to

$$\text{transport}^{X \mapsto X}(\text{ua}(e), m(\text{transport}^{X \mapsto X}(\text{ua}(e)^{-1}, b_1), \text{transport}^{X \mapsto X}(\text{ua}(e)^{-1}, b_2))).$$

Then, because ua is quasi-inverse to $\text{transport}^{X \mapsto X}$, this is equal to

$$e(m(e^{-1}(b_1), e^{-1}(b_2))).$$

Thus, given two elements of B , the induced multiplication m' sends them to A using the equivalence e , multiplies them in A , and then brings the result back to B by e , just as one would expect.

Moreover, though we do not show the proof, one can calculate that the induced proof that m' is associative (see (2.14.2)) is equal to a function sending $b_1, b_2, b_3 : B$ to a path given by the

following steps:

$$\begin{aligned}
m'(m'(b_1, b_2), b_3) &= e(m(e^{-1}(m'(b_1, b_2)), e^{-1}(b_3))) \\
&= e(m(e^{-1}(e(m(e^{-1}(b_1), e^{-1}(b_2)))), e^{-1}(b_3))) \\
&= e(m(m(e^{-1}(b_1), e^{-1}(b_2)), e^{-1}(b_3))) \\
&= e(m(e^{-1}(b_1), m(e^{-1}(b_2), e^{-1}(b_3)))) \\
&= e(m(e^{-1}(b_1), e^{-1}(e(m(e^{-1}(b_2), e^{-1}(b_3))))) \\
&= e(m(e^{-1}(b_1), e^{-1}(m'(b_2, b_3)))) \\
&= m'(b_1, m'(b_2, b_3)).
\end{aligned} \tag{2.14.3}$$

These steps use the proof a that m is associative and the inverse laws for e . From an algebra perspective, it may seem strange to investigate the identity of a proof that an operation is associative, but this makes sense if we think of A and B as general spaces, with non-trivial homotopies between paths. In Chapter 3, we will introduce the notion of a *set*, which is a type with only trivial homotopies, and if we consider semigroup structures on sets, then any two such associativity proofs are automatically equal.

2.14.2 Equality of semigroups

Using the equations for path spaces discussed in the previous sections, we can investigate when two semigroups are equal. Given semigroups (A, m, a) and (B, m', a') , by Theorem 2.7.2, the type of paths $(A, m, a) =_{\text{Semigroup}} (B, m', a')$ is equal to the type of pairs

$$\begin{aligned}
p_1 : A &=_U B \quad \text{and} \\
p_2 : \text{transport}^{\text{SemigroupStr}}(p_1, (m, a)) &= (m', a').
\end{aligned}$$

By univalence, p_1 is $\text{ua}(e)$ for some equivalence e . By Theorem 2.7.2, function extensionality, and the above analysis of transport in the type family SemigroupStr , p_2 is equivalent to a pair of proofs, the first of which shows that

$$\prod_{y_1, y_2 : B} e(m(e^{-1}(y_1), e^{-1}(y_2))) = m'(y_1, y_2)$$

and the second of which shows that a' is equal to the induced associativity proof constructed from a in (2.14.3). But by cancellation of inverses (2.14.2) is equivalent to

$$\prod_{x_1, x_2 : A} e(m(x_1, x_2)) = m'(e(x_1), e(x_2)).$$

This says that e commutes with the binary operation, in the sense that it takes multiplication in A (i.e. m) to multiplication in B (i.e. m'). A similar rearrangement is possible for the equation relating a and a' . Thus, an equality of semigroups consists exactly of an equivalence on the carrier types that commutes with the semigroup structure.

For general types, the proof of associativity is thought of as part of the structure of a semigroup. However, if we restrict to set-like types (again, see Chapter 3), the equation relating a and a' is trivially true. Moreover, in this case, an equivalence between sets is exactly a bijection. Thus, we have arrived at a standard definition of a *semigroup isomorphism*: a bijection on the carrier sets that preserves the multiplication operation. It is also possible to use the category-theoretic definition of isomorphism, by defining a *semigroup homomorphism* to be a map that preserves the

multiplication, and arrive at the conclusion that equality of semigroups is the same as two mutually inverse homomorphisms; but we will not show the details here; see ??.

The conclusion is that, thanks to univalence, semigroups are equal precisely when they are isomorphic as algebraic structures. As we will see in ??, the conclusion applies more generally: in homotopy type theory, all constructions of mathematical structures automatically respect isomorphisms, without any tedious proofs or abuse of notation.

2.15 Universal properties

By combining the path computation rules described in the preceding sections, we can show that various type forming operations satisfy the expected universal properties, interpreted in a homotopical way as equivalences. For instance, given types X, A, B , we have a function

$$(X \rightarrow A \times B) \rightarrow (X \rightarrow A) \times (X \rightarrow B) \quad (2.15.1)$$

defined by $f \mapsto (\text{pr}_1 \circ f, \text{pr}_2 \circ f)$.

Theorem 2.15.2. (2.15.1) is an equivalence.

Proof. We define the quasi-inverse by sending (g, h) to $\lambda x. (g(x), h(x))$. (Technically, we have used the induction principle for the cartesian product $(X \rightarrow A) \times (X \rightarrow B)$, to reduce to the case of a pair. From now on we will often apply this principle without explicit mention.)

Now given $f : X \rightarrow A \times B$, the round-trip composite yields the function

$$\lambda x. (\text{pr}_1(f(x)), \text{pr}_2(f(x))). \quad (2.15.3)$$

By Theorem 2.6.2, for any $x : X$ we have $(\text{pr}_1(f(x)), \text{pr}_2(f(x))) = f(x)$. Thus, by function extensionality, the function (2.15.3) is equal to f .

On the other hand, given (g, h) , the round-trip composite yields the pair $(\lambda x. g(x), \lambda x. h(x))$. By the uniqueness principle for functions, this is (judgmentally) equal to (g, h) . \square

In fact, we also have a dependently typed version of this universal property. Suppose given a type X and type families $A, B : X \rightarrow \mathcal{U}$. Then we have a function

$$\left(\prod_{x:X} (A(x) \times B(x)) \right) \rightarrow \left(\prod_{x:X} A(x) \right) \times \left(\prod_{x:X} B(x) \right) \quad (2.15.4)$$

defined as before by $f \mapsto (\text{pr}_1 \circ f, \text{pr}_2 \circ f)$.

Theorem 2.15.5. (2.15.4) is an equivalence.

Proof. Left to the reader. \square

Just as Σ -types are a generalization of cartesian products, they satisfy a generalized version of this universal property. Jumping right to the dependently typed version, suppose we have a type X and type families $A : X \rightarrow \mathcal{U}$ and $P : \prod_{(x:X)} A(x) \rightarrow \mathcal{U}$. Then we have a function

$$\left(\prod_{x:X} \sum_{(a:A(x))} P(x, a) \right) \rightarrow \left(\sum_{(g:\prod_{(x:X)} A(x))} \prod_{(x:X)} P(x, g(x)) \right). \quad (2.15.6)$$

Note that if we have $P(x, a) := B(x)$ for some $B : X \rightarrow \mathcal{U}$, then (2.15.6) reduces to (2.15.4).

Theorem 2.15.7. (2.15.6) is an equivalence.

Proof. As before, we define a quasi-inverse to send (g, h) to the function $\lambda x. (g(x), h(x))$. Now given $f : \prod_{(x:X)} \sum_{(a:A(x))} P(x, a)$, the round-trip composite yields the function

$$\lambda x. (\text{pr}_1(f(x)), \text{pr}_2(f(x))). \quad (2.15.8)$$

Now for any $x : X$, by Corollary 2.7.3 (the uniqueness principle for Σ -types) we have

$$(\text{pr}_1(f(x)), \text{pr}_2(f(x))) = f(x).$$

Thus, by function extensionality, (2.15.8) is equal to f . On the other hand, given (g, h) , the round-trip composite yields $(\lambda x. g(x), \lambda x. h(x))$, which is judgmentally equal to (g, h) as before. \square

This is noteworthy because the propositions-as-types interpretation of (2.15.6) is “the axiom of choice”. If we read Σ as “there exists” and Π (sometimes) as “for all”, we can pronounce:

- $\prod_{(x:X)} \sum_{(a:A(x))} P(x, a)$ as “for all $x : X$ there exists an $a : A(x)$ such that $P(x, a)$ ”, and
- $\sum_{(g:\prod_{(x:X)} A(x))} \prod_{(x:X)} P(x, g(x))$ as “there exists a choice function $g : \prod_{(x:X)} A(x)$ such that for all $x : X$ we have $P(x, g(x))$ ”.

Thus, Theorem 2.15.7 says that not only is the axiom of choice “true”, its antecedent is actually equivalent to its conclusion. (On the other hand, the classical mathematician may find that (2.15.6) does not carry the usual meaning of the axiom of choice, since we have already specified the values of g , and there are no choices left to be made. We will return to this point in §3.8.)

The above universal property for pair types is for “mapping in”, which is familiar from the category-theoretic notion of products. However, pair types also have a universal property for “mapping out”, which may look less familiar. In the case of cartesian products, the non-dependent version simply expresses the cartesian closure adjunction:

$$((A \times B) \rightarrow C) \simeq (A \rightarrow (B \rightarrow C)).$$

The dependent version of this is formulated for a type family $C : A \times B \rightarrow \mathcal{U}$:

$$\left(\prod_{w:A \times B} C(w) \right) \simeq \left(\prod_{(x:A)} \prod_{(y:B)} C(x, y) \right).$$

Here the right-to-left function is simply the induction principle for $A \times B$, while the left-to-right is evaluation at a pair. We leave it to the reader to prove that these are quasi-inverses. There is also a version for Σ -types:

$$\left(\prod_{w:\sum_{(x:A)} B(x)} C(w) \right) \simeq \left(\prod_{(x:A)} \prod_{(y:B(x))} C(x, y) \right). \quad (2.15.9)$$

Again, the right-to-left function is the induction principle.

Some other induction principles are also part of universal properties of this sort. For instance, path induction is the right-to-left direction of an equivalence as follows:

$$\left(\prod_{(x:A)} \prod_{(p:a=x)} B(x, p) \right) \simeq B(a, \text{refl}_a) \quad (2.15.10)$$

for any $a : A$ and type family $B : \prod_{(x:A)}(a = x) \rightarrow \mathcal{U}$. However, inductive types with recursion, such as the natural numbers, have more complicated universal properties; see Chapter 5.

Since Theorem 2.15.2 expresses the usual universal property of a cartesian product (in an appropriate homotopy-theoretic sense), the categorically inclined reader may well wonder about other limits and colimits of types. In Exercise 2.9 we ask the reader to show that the coproduct type $A + B$ also has the expected universal property, and the nullary cases of $\mathbf{1}$ (the terminal object) and $\mathbf{0}$ (the initial object) are easy.

For pullbacks, the expected explicit construction works: given $f : A \rightarrow C$ and $g : B \rightarrow C$, we define

$$A \times_C B := \sum_{(a:A)} \sum_{(b:B)} (f(a) = g(b)). \quad (2.15.11)$$

In Exercise 2.11 we ask the reader to verify this. Some more general homotopy limits can be constructed in a similar way, but for colimits we will need a new ingredient; see Chapter 6.

Notes

The definition of identity types, with their induction principle, is due to Martin-Löf [ML98]. As mentioned in the notes to Chapter 1, our identity types are those that belong to *intensional* type theory, rather than *extensional* type theory. In general, a notion of equality is said to be “intensional” if it distinguishes objects based on their particular definitions, and “extensional” if it does not distinguish between objects that have the same “extension” or “observable behavior”. In the terminology of Frege, an intensional equality compares *sense*, while an extensional one compares only *reference*. We may also speak of one equality being “more” or “less” extensional than another, meaning that it takes account of fewer or more intensional aspects of objects, respectively.

Intensional type theory is so named because its *judgmental* equality, $x \equiv y$, is a very intensional equality: it says essentially that x and y “have the same definition”, after we expand the defining equations of functions. By contrast, the propositional equality type $x = y$ is more extensional, even in the axiom-free intensional type theory of Chapter 1: for instance, we can prove by induction that $n + m = m + n$ for all $m, n : \mathbb{N}$, but we cannot say that $n + m \equiv m + n$ for all $m, n : \mathbb{N}$, since the *definition* of addition treats its arguments asymmetrically. We can make the identity type of intensional type theory even more extensional by adding axioms such as function extensionality (two functions are equal if they have the same behavior on all inputs, regardless of how they are defined) and univalence (which can be regarded as an extensionality property for the universe: two types are equal if they behave the same in all contexts). The axioms of function extensionality, and univalence in the special case of mere propositions (“propositional extensionality”), appeared already in the first type theories of Russell and Church.

As mentioned before, *extensional* type theory includes also a “reflection rule” saying that if $p : x = y$, then in fact $x \equiv y$. Thus extensional type theory is so named because it does *not* admit any purely *intensional* equality: the reflection rule forces the judgmental equality to coincide with the more extensional identity type. Moreover, from the reflection rule one may deduce function extensionality (at least in the presence of a judgmental uniqueness principle for functions). However, the reflection rule also implies that all the higher groupoid structure collapses (see Exercise 2.14), and hence is inconsistent with the univalence axiom (see Example 3.1.9). Therefore, regarding univalence as an extensionality property, one may say that intensional type theory permits identity types that are “more extensional” than extensional type theory does.

The proofs of symmetry (inversion) and transitivity (concatenation) for equalities are well-known in type theory. The fact that these make each type into a 1-groupoid (up to homotopy) was exploited in [HS98] to give the first “homotopy” style semantics for type theory.

The actual homotopical interpretation, with identity types as path spaces, and type families as fibrations, is due to [AW09], who used the formalism of Quillen model categories. An interpretation in (strict) ∞ -groupoids was also given in the thesis [War08]. For a construction of *all* the higher operations and coherences of an ∞ -groupoid in type theory, see [Lum10] and [vdBG11].

Operations such as $\text{transport}^P(p, -)$ and ap_f , and one good notion of equivalence, were first studied extensively in type theory by Voevodsky, using the proof assistant Coq. Subsequently, many other equivalent definitions of equivalence have been found, which are compared in Chapter 4.

The “computational” interpretation of identity types, transport, and so on described in §2.5 has been emphasized by [LH12]. They also described a “1-truncated” type theory (see Chapter 7) in which these rules are judgmental equalities. The possibility of extending this to the full untruncated theory is a subject of current research.

The naive form of function extensionality which says that “if two functions are pointwise equal, then they are equal” is a common axiom in type theory, going all the way back to [WR27]. Some stronger forms of function extensionality were considered in [Gar09]. The version we have used, which identifies the identity types of function types up to equivalence, was first studied by Voevodsky, who also proved that it is implied by the naive version (and by univalence; see §4.9).

The univalence axiom is also due to Voevodsky. It was originally motivated by semantic considerations in the simplicial set model; see [KLV12]. A similar axiom motivated by the groupoid model was proposed by Hofmann and Streicher [HS98] under the name “universe extensionality”. It used quasi-inverses (2.4.5) rather than a good notion of “equivalence”, and hence is correct (and equivalent to univalence) only for a universe of 1-types (see Definition 3.1.7).

In the type theory we are using in this book, function extensionality and univalence have to be assumed as axioms, i.e. elements asserted to belong to some type but not constructed according to the rules for that type. While serviceable, this has a few drawbacks. For instance, type theory is formally better-behaved if we can base it entirely on rules rather than asserting axioms. It is also sometimes inconvenient that the theorems of §§2.6–2.13 are only propositional equalities (paths) or equivalences, since then we must explicitly mention whenever we pass back and forth across them. One direction of current research in homotopy type theory is to describe a type system in which these rules are *judgmental* equalities, solving both of these problems at once. So far this has only been done in some simple cases, although preliminary results such as [LH12] are promising. There are also other potential ways to introduce univalence and function extensionality into a type theory, such as having a sufficiently powerful notion of “higher quotients” or “higher inductive-recursive types”.

The simple conclusions in §§2.12–2.13 such as “`inl` and `inr` are injective and disjoint” are well-known in type theory, and the construction of the function `encode` is the usual way to prove them. The more refined approach we have described, which characterizes the entire identity type of a positive type (up to equivalence), is a more recent development; see e.g. [LS13].

The type-theoretic axiom of choice (2.15.6) was noticed in William Howard’s original paper [How80] on the propositions-as-types correspondence, and was studied further by Martin-Löf with the introduction of his dependent type theory. It is mentioned as a “distributivity law” in Bourbaki’s set theory [Bou68].

For a more comprehensive (and formalized) discussion of pullbacks and more general ho-

motopy limits in homotopy type theory, see [AKL13]. Limits of diagrams over directed graphs are the easiest general sort of limit to formalize; the problem with diagrams over categories (or more generally $(\infty, 1)$ -categories) is that in general, infinitely many coherence conditions are involved in the notion of (homotopy coherent) diagram. Resolving this problem is an important open question in homotopy type theory.

Exercises

Exercise 2.1. Show that the three obvious proofs of Lemma 2.1.2 are pairwise equal.

Exercise 2.2. Show that the three equalities of proofs constructed in the previous exercise form a commutative triangle. In other words, if the three definitions of concatenation are denoted by $(p \bullet_1 q)$, $(p \bullet_2 q)$, and $(p \bullet_3 q)$, then the concatenated equality

$$(p \bullet_1 q) = (p \bullet_2 q) = (p \bullet_3 q)$$

is equal to the equality $(p \bullet_1 q) = (p \bullet_3 q)$.

Exercise 2.3. Give a fourth, different, proof of Lemma 2.1.2, and prove that it is equal to the others.

Exercise 2.4. Define, by induction on n , a general notion of **n -dimensional path** in a type A , simultaneously with the type of boundaries for such paths.

Exercise 2.5. Prove that the functions (2.3.6) and (2.3.7) are inverse equivalences.

Exercise 2.6. Prove that if $p : x = y$, then the function $(p \bullet -) : (y = z) \rightarrow (x = z)$ is an equivalence.

Exercise 2.7. State and prove a generalization of Theorem 2.6.5 from cartesian products to Σ -types.

Exercise 2.8. State and prove an analogue of Theorem 2.6.5 for coproducts.

Exercise 2.9. Prove that coproducts have the expected universal property,

$$(A + B \rightarrow X) \simeq (A \rightarrow X) \times (B \rightarrow X).$$

Can you generalize this to an equivalence involving dependent functions?

Exercise 2.10. Prove that Σ -types are “associative”, in that for any $A : \mathcal{U}$ and families $B : A \rightarrow \mathcal{U}$ and $C : (\sum_{(x:A)} B(x)) \rightarrow \mathcal{U}$, we have

$$\left(\sum_{(x:A)} \sum_{(y:B(x))} C((x,y)) \right) \simeq \left(\sum_{p:\sum_{(x:A)} B(x)} C(p) \right).$$

Exercise 2.11. A (homotopy) **commutative square**

$$\begin{array}{ccc} P & \xrightarrow{h} & A \\ k \downarrow & & \downarrow f \\ B & \xrightarrow{g} & C \end{array}$$

consists of functions f , g , h , and k as shown, together with a path $f \circ h = g \circ k$. Note that this is exactly an element of the pullback $(P \rightarrow A) \times_{P \rightarrow C} (P \rightarrow B)$ as defined in (2.15.11). A commutative square is called a (homotopy) **pullback square** if for any X , the induced map

$$(X \rightarrow P) \rightarrow (X \rightarrow A) \times_{(X \rightarrow C)} (X \rightarrow B)$$

is an equivalence. Prove that the pullback $P := A \times_C B$ defined in (2.15.11) is the corner of a pullback square.

Exercise 2.12. Suppose given two commutative squares

$$\begin{array}{ccccc} A & \longrightarrow & C & \longrightarrow & E \\ \downarrow & & \downarrow & & \downarrow \\ B & \longrightarrow & D & \longrightarrow & F \end{array}$$

and suppose that the right-hand square is a pullback square. Prove that the left-hand square is a pullback square if and only if the outer rectangle is a pullback square.

Exercise 2.13. Show that $(\mathbf{2} \simeq \mathbf{2}) \simeq \mathbf{2}$.

Exercise 2.14. Suppose we add to type theory the *equality reflection rule* which says that if there is an element $p : x = y$, then in fact $x \equiv y$. Prove that for any $p : x = x$ we have $p \equiv \text{refl}_x$. (This implies that every type is a *set* in the sense to be introduced in §3.1; see §7.2.)

Exercise 2.15. Show that Lemma 2.10.5 can be strengthened to

$$\text{transport}^B(p, -) =_{B(x) \rightarrow B(y)} \text{idtoeqv}(\text{ap}_B(p))$$

without using function extensionality. (In this and other similar cases, the apparently weaker formulation has been chosen for readability and consistency.)

Exercise 2.16. Suppose that rather than function extensionality (Axiom 2.9.3), we suppose only the existence of an element

$$\text{funext} : \prod_{(A:\mathcal{U})} \prod_{(B:A \rightarrow \mathcal{U})} \prod_{(f,g:\prod_{(x:A)} B(x))} (f \sim g) \rightarrow (f = g)$$

(with no relationship to `happly` assumed). Prove that in fact, this is sufficient to imply the whole function extensionality axiom (that `happly` is an equivalence). This is due to Voevodsky; its proof is tricky and may require concepts from later chapters.

Exercise 2.17.

- (i) Show that if $A \simeq A'$ and $B \simeq B'$, then $(A \times B) \simeq (A' \times B')$.
- (ii) Give two proofs of this fact, one using univalence and one not using it, and show that the two proofs are equal.
- (iii) Formulate and prove analogous results for the other type formers: Σ , \rightarrow , Π , and $+$.

Exercise 2.18. State and prove a version of Lemma 2.4.3 for dependent functions.

Chapter 3

Sets and logic

Type theory, formal or informal, is a collection of rules for manipulating types and their elements. But when writing mathematics informally in natural language, we generally use familiar words, particularly logical connectives such as “and” and “or”, and logical quantifiers such as “for all” and “there exists”. In contrast to set theory, type theory offers us more than one way to regard these English phrases as operations on types. This potential ambiguity needs to be resolved, by setting out local or global conventions, by introducing new annotations to informal mathematics, or both. This requires some getting used to, but is offset by the fact that because type theory permits this finer analysis of logic, we can represent mathematics more faithfully, with fewer “abuses of language” than in set-theoretic foundations. In this chapter we will explain the issues involved, and justify the choices we have made.

3.1 Sets and n -types

In order to explain the connection between the logic of type theory and the logic of set theory, it is helpful to have a notion of *set* in type theory. While types in general behave like spaces or higher groupoids, there is a subclass of them that behave more like the sets in a traditional set-theoretic system. Categorically, we may consider *discrete* groupoids, which are determined by a set of objects and only identity morphisms as higher morphisms; while topologically, we may consider spaces having the discrete topology. More generally, we may consider groupoids or spaces that are *equivalent* to ones of this sort; since everything we do in type theory is up to homotopy, we can’t expect to tell the difference.

Intuitively, we would expect a type to “be a set” in this sense if it has no higher homotopical information: any two parallel paths are equal (up to homotopy), and similarly for parallel higher paths at all dimensions. Fortunately, because everything in homotopy type theory is automatically functorial/continuous, it turns out to be sufficient to ask this at the bottom level.

Definition 3.1.1. A type A is a **set** if for all $x, y : A$ and all $p, q : x = y$, we have $p = q$.

More precisely, the proposition $\text{isSet}(A)$ is defined to be the type

$$\text{isSet}(A) := \prod_{(x,y:A)} \prod_{(p,q:x=y)} (p = q).$$

As mentioned in §1.1, the sets in homotopy type theory are not like the sets in ZF set theory, in that there is no global “membership predicate” \in . They are more like the sets used in structural

mathematics and in category theory, whose elements are “abstract points” to which we give structure with functions and relations. This is all we need in order to use them as a foundational system for most set-based mathematics; we will see some examples in ??.

Which types are sets? In Chapter 7 we will study a more general form of this question in depth, but for now we can observe some easy examples.

Example 3.1.2. The type **1** is a set. For by Theorem 2.8.1, for any $x, y : \mathbf{1}$ the type $(x = y)$ is equivalent to **1**. Since any two elements of **1** are equal, this implies that any two elements of $x = y$ are equal.

Example 3.1.3. The type **0** is a set, for given any $x, y : \mathbf{0}$ we may deduce anything we like, by the induction principle of **0**.

Example 3.1.4. The type \mathbb{N} of natural numbers is also a set. This follows from Theorem 2.13.1, since all equality types $x =_{\mathbb{N}} y$ are equivalent to either **1** or **0**, and any two inhabitants of **1** or **0** are equal. We will see another proof of this fact in Chapter 7.

Most of the type forming operations we have considered so far also preserve sets.

Example 3.1.5. If A and B are sets, then so is $A \times B$. For given $x, y : A \times B$ and $p, q : x = y$, by Theorem 2.6.2 we have $p = \text{pair}^=(\text{ap}_{\text{pr}_1}(p), \text{ap}_{\text{pr}_2}(p))$ and $q = \text{pair}^=(\text{ap}_{\text{pr}_1}(q), \text{ap}_{\text{pr}_2}(q))$. But $\text{ap}_{\text{pr}_1}(p) = \text{ap}_{\text{pr}_1}(q)$ since A is a set, and $\text{ap}_{\text{pr}_2}(p) = \text{ap}_{\text{pr}_2}(q)$ since B is a set; hence $p = q$.

Similarly, if A is a set and $B : A \rightarrow \mathcal{U}$ is such that each $B(x)$ is a set, then $\sum_{(x:A)} B(x)$ is a set.

Example 3.1.6. If A is any type and $B : A \rightarrow \mathcal{U}$ is such that each $B(x)$ is a set, then the type $\prod_{(x:A)} B(x)$ is a set. For suppose $f, g : \prod_{(x:A)} B(x)$ and $p, q : f = g$. By function extensionality, we have

$$p = \text{funext}(x \mapsto \text{happly}(p, x)) \quad \text{and} \quad q = \text{funext}(x \mapsto \text{happly}(q, x)).$$

But for any $x : A$, we have

$$\text{happly}(p, x) : f(x) = g(x) \quad \text{and} \quad \text{happly}(q, x) : f(x) = g(x),$$

so since $B(x)$ is a set we have $\text{happly}(p, x) = \text{happly}(q, x)$. Now using function extensionality again, the dependent functions $(x \mapsto \text{happly}(p, x))$ and $(x \mapsto \text{happly}(q, x))$ are equal, and hence (applying $\text{ap}_{\text{funext}}$) so are p and q .

For more examples, see Exercises 3.2 and 3.3. For a more systematic investigation of the subsystem (category) of all sets in homotopy type theory, see ??.

Sets are just the first rung on a ladder of what are called *homotopy n-types*. The next rung consists of *1-types*, which are analogous to 1-groupoids in category theory. The defining property of a set (which we may also call a *0-type*) is that it has no non-trivial paths. Similarly, the defining property of a 1-type is that it has no non-trivial paths between paths:

Definition 3.1.7. A type A is a **1-type** if for all $x, y : A$ and $p, q : x = y$ and $r, s : p = q$, we have $r = s$.

Similarly, we can define 2-types, 3-types, and so on. We will define the general notion of n -type inductively in Chapter 7, and study the relationships between n -types for different values of n .

However, for now it is useful to have two facts in mind. First, the levels are upward-closed: if A is an n -type then A is an $(n + 1)$ -type. For example:

Lemma 3.1.8. If A is a set (that is, $\text{isSet}(A)$ is inhabited), then A is a 1-type.

Proof. Suppose $f : \text{isSet}(A)$; then for any $x, y : A$ and $p, q : x = y$ we have $f(x, y, p, q) : p = q$. Fix x, y , and p , and define $g : \prod_{(q:x=y)}(p = q)$ by $g(q) := f(x, y, p, q)$. Then for any $r : q = q'$, we have $\text{apd}_g(r) : r_*(g(q)) = g(q')$. By Lemma 2.11.2, therefore, we have $g(q) \cdot r = g(q')$.

In particular, suppose given x, y, p, q and $r, s : p = q$, as in Definition 3.1.7, and define g as above. Then $g(p) \cdot r = g(q)$ and also $g(p) \cdot s = g(q)$, hence by cancellation $r = s$. \square

Second, this stratification of types by level is not degenerate, in the sense that not all types are sets:

Example 3.1.9. The universe \mathcal{U} is not a set. To prove this, it suffices to exhibit a type A and a path $p : A = A$ which is not equal to refl_A . Take $A = \mathbf{2}$, and let $f : A \rightarrow A$ be defined by $f(0_2) := 1_2$ and $f(1_2) := 0_2$. Then $f(f(x)) = x$ for all x (by an easy case analysis), so f is an equivalence. Hence, by univalence, f gives rise to a path $p : A = A$.

If p were equal to refl_A , then (again by univalence) f would equal the identity function of A . But this would imply that $0_2 = 1_2$, contradicting Remark 2.12.6.

In Chapter 6 and ?? we will show that for any n , there are types which are not n -types.

Note that A is a 1-type exactly when for any $x, y : A$, the identity type $x =_A y$ is a set. (Thus, Lemma 3.1.8 could equivalently be read as saying that the identity types of a set are also sets.) This will be the basis of the recursive definition of n -types we will give in Chapter 7.

We can also extend this characterization “downwards” from sets. That is, a type A is a set just when for any $x, y : A$, any two elements of $x =_A y$ are equal. Since sets are equivalently 0-types, it is natural to call a type a *(−1)-type* if it has this latter property (any two elements of it are equal). Such types may be regarded as *propositions in a narrow sense*, and their study is just what is usually called “logic”; it will occupy us for the rest of this chapter.

3.2 Propositions as types?

Until now, we have been following the straightforward “propositions as types” philosophy described in §1.11, according to which English phrases such as “there exists an $x : A$ such that $P(x)$ ” are interpreted by corresponding types such as $\sum_{(x:A)} P(x)$, with the proof of a statement being regarded as judging some specific element to inhabit that type. However, we have also seen some ways in which the “logic” resulting from this reading seems unfamiliar to a classical mathematician. For instance, in Theorem 2.15.7 we saw that the statement

“If for all $x : X$ there exists an $a : A(x)$ such that $P(x, a)$, then there exists a function $g : \prod_{(x:X)} A(x)$ such that for all $x : X$ we have $P(x, g(x))$ ”, (3.2.1)

which looks like the classical *axiom of choice*, is always true under this reading. This is a noteworthy, and often useful, feature of the propositions-as-types logic, but it also illustrates how significantly it differs from the classical interpretation of logic, under which the axiom of choice is not a logical truth, but an additional “axiom”.

On the other hand, we can now also show that corresponding statements looking like the classical *law of double negation* and *law of excluded middle* are incompatible with the univalence axiom.

Theorem 3.2.2. *It is not the case that for all $A : \mathcal{U}$ we have $\neg(\neg A) \rightarrow A$.*

Proof. Recall that $\neg A \equiv (A \rightarrow \mathbf{0})$. We also read “it is not the case that ...” as the operator \neg . Thus, in order to prove this statement, it suffices to assume given some $f : \prod_{(A:\mathcal{U})} (\neg\neg A \rightarrow A)$ and construct an element of $\mathbf{0}$.

The idea of the following proof is to observe that f , like any function in type theory, is “continuous”. By univalence, this implies that f is *natural* with respect to equivalences of types. From this, and a fixed-point-free autoequivalence, we will be able to extract a contradiction.

Let $e : \mathbf{2} \simeq \mathbf{2}$ be the equivalence defined by $e(1_2) := 0_2$ and $e(0_2) := 1_2$, as in Example 3.1.9. Let $p : \mathbf{2} = \mathbf{2}$ be the path corresponding to e by univalence, i.e. $p := \text{ua}(e)$. Then we have $f(\mathbf{2}) : \neg\neg\mathbf{2} \rightarrow \mathbf{2}$ and

$$\text{apd}_f(p) : \text{transport}^{A \mapsto (\neg\neg A \rightarrow A)}(p, f(\mathbf{2})) = f(\mathbf{2}).$$

Hence, for any $u : \neg\neg\mathbf{2}$, we have

$$\text{happly}(\text{apd}_f(p), u) : \text{transport}^{A \mapsto (\neg\neg A \rightarrow A)}(p, f(\mathbf{2}))(u) = f(\mathbf{2})(u).$$

Now by (2.9.4), transporting $f(\mathbf{2}) : \neg\neg\mathbf{2} \rightarrow \mathbf{2}$ along p in the type family $A \mapsto (\neg\neg A \rightarrow A)$ is equal to the function which transports its argument along p^{-1} in the type family $A \mapsto \neg\neg A$, applies $f(\mathbf{2})$, then transports the result along p in the type family $A \mapsto A$:

$$\text{transport}^{A \mapsto (\neg\neg A \rightarrow A)}(p, f(\mathbf{2}))(u) = \text{transport}^{A \mapsto A}(p, f(\mathbf{2})(\text{transport}^{A \mapsto \neg\neg A}(p^{-1}, u))).$$

However, any two points $u, v : \neg\neg\mathbf{2}$ are equal by function extensionality, since for any $x : \neg\mathbf{2}$ we have $u(x) : \mathbf{0}$ and thus we can derive any conclusion, in particular $u(x) = v(x)$. Thus, we have $\text{transport}^{A \mapsto \neg\neg A}(p^{-1}, u) = u$, and so from $\text{happly}(\text{apd}_f(p), u)$ we obtain an equality

$$\text{transport}^{A \mapsto A}(p, f(\mathbf{2}))(u) = f(\mathbf{2})(u).$$

Finally, as discussed in §2.10, transporting in the type family $A \mapsto A$ along the path $p \equiv \text{ua}(e)$ is equivalent to applying the equivalence e ; thus we have

$$e(f(\mathbf{2})(u)) = f(\mathbf{2})(u). \tag{3.2.3}$$

However, we can also prove that

$$\prod_{x:\mathbf{2}} \neg(e(x) = x). \tag{3.2.4}$$

This follows from a case analysis on x : both cases are immediate from the definition of e and the fact that $0_2 \neq 1_2$ (Remark 2.12.6). Thus, applying (3.2.4) to $f(\mathbf{2})(u)$ and (3.2.3), we obtain an element of $\mathbf{0}$. \square

Remark 3.2.5. In particular, this implies that there can be no Hilbert-style “choice operator” which selects an element of every nonempty type. The point is that no such operator can be *natural*, and under the univalence axiom, all functions acting on types must be natural with respect to equivalences.

Remark 3.2.6. It is, however, still the case that $\neg\neg\neg A \rightarrow \neg A$ for any A ; see Exercise 1.11.

Corollary 3.2.7. *It is not the case that for all $A : \mathcal{U}$ we have $A + (\neg A)$.*

Proof. Suppose we had $g : \prod_{(A:\mathcal{U})}(A + (\neg A))$. We will show that then $\prod_{(A:\mathcal{U})}(\neg\neg A \rightarrow A)$, so that we can apply Theorem 3.2.2. Thus, suppose $A : \mathcal{U}$ and $u : \neg\neg A$; we want to construct an element of A .

Now $g(A) : A + (\neg A)$, so by case analysis, we may assume either $g(A) \equiv \text{inl}(a)$ for some $a : A$, or $g(A) \equiv \text{inr}(w)$ for some $w : \neg A$. In the first case, we have $a : A$, while in the second case we have $u(w) : \mathbf{0}$ and so we can obtain anything we wish (such as A). Thus, in both cases we have an element of A , as desired. \square

Thus, if we want to assume the univalence axiom (which, of course, we do) and still leave ourselves the option of classical reasoning (which is also desirable), we cannot use the unmodified propositions-as-types principle to interpret *all* informal mathematical statements into type theory, since then the law of excluded middle would be false. However, neither do we want to discard propositions-as-types entirely, because of its many good properties (such as simplicity, constructivity, and computability). We now discuss a modification of propositions-as-types which resolves these problems; in §3.10 we will return to the question of which logic to use when.

3.3 Mere propositions

We have seen that the propositions-as-types logic has both good and bad properties. Both have a common cause: when types are viewed as propositions, they can contain more information than mere truth or falsity, and all “logical” constructions on them must respect this additional information. This suggests that we could obtain a more conventional logic by restricting attention to types that do *not* contain any more information than a truth value, and only regarding these as logical propositions.

Such a type A will be “true” if it is inhabited, and “false” if its inhabitation yields a contradiction (i.e. if $\neg A \equiv (A \rightarrow \mathbf{0})$ is inhabited). What we want to avoid, in order to obtain a more traditional sort of logic, is treating as logical propositions those types for which giving an element of them gives more information than simply knowing that the type is inhabited. For instance, if we are given an element of $\mathbf{2}$, then we receive more information than the mere fact that $\mathbf{2}$ contains some element. Indeed, we receive exactly *one bit* more information: we know *which* element of $\mathbf{2}$ we were given. By contrast, if we are given an element of $\mathbf{1}$, then we receive no more information than the mere fact that $\mathbf{1}$ contains an element, since any two elements of $\mathbf{1}$ are equal to each other. This suggests the following definition.

Definition 3.3.1. A type P is a **mere proposition** if for all $x, y : P$ we have $x = y$.

Note that since we are still doing mathematics *in* type theory, this is a definition *in* type theory, which means it is a type — or, rather, a type family. Specifically, for any $P : \mathcal{U}$, the type $\text{isProp}(P)$ is defined to be

$$\text{isProp}(P) := \prod_{x,y:P} (x = y).$$

Thus, to assert that “ P is a mere proposition” means to exhibit an inhabitant of $\text{isProp}(P)$, which is a dependent function connecting any two elements of P by a path. The continuity/naturality of this function implies that not only are any two elements of P equal, but P contains no higher homotopy either.

Lemma 3.3.2. If P is a mere proposition and $x_0 : P$, then $P \simeq \mathbf{1}$.

Proof. Define $f : P \rightarrow \mathbf{1}$ by $f(x) := \star$, and $g : \mathbf{1} \rightarrow P$ by $g(\star) := x_0$. The claim follows from the next lemma, and the observation that $\mathbf{1}$ is a mere proposition by Theorem 2.8.1. \square

Lemma 3.3.3. If P and Q are mere propositions such that $P \rightarrow Q$ and $Q \rightarrow P$, then $P \simeq Q$.

Proof. Suppose given $f : P \rightarrow Q$ and $g : Q \rightarrow P$. Then for any $x : P$, we have $g(f(x)) = x$ since P is a mere proposition. Similarly, for any $y : Q$ we have $f(g(y)) = y$ since Q is a mere proposition; thus f and g are quasi-inverses. \square

That is, as promised in §1.11, if two mere propositions are logically equivalent, then they are equivalent.

In homotopy theory, a space that is homotopy equivalent to $\mathbf{1}$ is said to be *contractible*. Thus, any mere proposition which is inhabited is contractible (see also §3.11). On the other hand, the uninhabited type $\mathbf{0}$ is also (vacuously) a mere proposition. In classical mathematics, at least, these are the only two possibilities.

Mere propositions are also called *subterminal objects* (if thinking categorically), *subsingletons* (if thinking set-theoretically), or *h-propositions*. The discussion in §3.1 suggests we should also call them (-1) -types; we will return to this in Chapter 7. The adjective “mere” emphasizes that although any type may be regarded as a proposition (which we prove by giving an inhabitant of it), a type that is a mere proposition cannot usefully be regarded as any *more* than a proposition: there is no additional information contained in a witness of its truth.

Note that a type A is a set if and only if for all $x, y : A$, the identity type $x =_A y$ is a mere proposition. On the other hand, by copying and simplifying the proof of Lemma 3.1.8, we have:

Lemma 3.3.4. *Every mere proposition is a set.*

Proof. Suppose $f : \text{isProp}(A)$; thus for all $x, y : A$ we have $f(x, y) : x = y$. Fix $x : A$ and define $g(y) := f(x, y)$. Then for any $y, z : A$ and $p : y = z$ we have $\text{apd}_g(p) : p_* g(y) = g(z)$. Hence by Lemma 2.11.2, we have $g(y) \cdot p = g(z)$, which is to say that $p = g(y)^{-1} \cdot g(z)$. Thus, for any $p, q : x = y$, we have $p = g(x)^{-1} \cdot g(y) = q$. \square

In particular, this implies:

Lemma 3.3.5. *For any type A , the types $\text{isProp}(A)$ and $\text{isSet}(A)$ are mere propositions.*

Proof. Suppose $f, g : \text{isProp}(A)$. By function extensionality, to show $f = g$ it suffices to show $f(x, y) = g(x, y)$ for any $x, y : A$. But $f(x, y)$ and $g(x, y)$ are both paths in A , and hence are equal because, by either f or g , we have that A is a mere proposition, and hence by Lemma 3.3.4 is a set. Similarly, suppose $f, g : \text{isSet}(A)$, which is to say that for all $a, b : A$ and $p, q : a = b$, we have $f(a, b, p, q) : p = q$ and $g(a, b, p, q) : p = q$. But by then since A is a set (by either f or g), and hence a 1-type, it follows that $f(a, b, p, q) = g(a, b, p, q)$; hence $f = g$ by function extensionality. \square

We have seen one other example so far: condition (iii) in §2.4 asserts that for any function f , the type $\text{isequiv}(f)$ should be a mere proposition.

3.4 Classical vs. intuitionistic logic

With the notion of mere proposition in hand, we can now give the proper formulation of the **law of excluded middle** in homotopy type theory:

$$\text{LEM} := \prod_{A:\mathcal{U}} (\text{isProp}(A) \rightarrow (A + \neg A)). \quad (3.4.1)$$

Similarly, the **law of double negation** is

$$\prod_{A:\mathcal{U}} (\text{isProp}(A) \rightarrow (\neg\neg A \rightarrow A)). \quad (3.4.2)$$

The two are also easily seen to be equivalent to each other—see Exercise 3.18—so from now on we will generally speak only of LEM.

This formulation of LEM avoids the “paradoxes” of Theorem 3.2.2 and Corollary 3.2.7, since **2** is not a mere proposition. In order to distinguish it from the more general propositions-as-types formulation, we rename the latter:

$$\text{LEM}_\infty := \prod_{A:\mathcal{U}} (A + \neg A).$$

For emphasis, the proper version (3.4.1) may be denoted LEM_{-1} ; see also Exercise 7.7. Although LEM is not a consequence of the basic type theory described in Chapter 1, it may be consistently assumed as an axiom (unlike its ∞ -counterpart). For instance, we will assume it in ??.

However, it can be surprising how far we can get without using LEM. Quite often, a simple reformulation of a definition or theorem enables us to avoid invoking excluded middle. While this takes a little getting used to sometimes, it is often worth the hassle, resulting in more elegant and more general proofs. We discussed some of the benefits of this in the introduction.

For instance, in classical mathematics, double negations are frequently used unnecessarily. A very simple example is the common assumption that a set A is “nonempty”, which literally means it is *not* the case that A contains *no* elements. Almost always what is really meant is the positive assertion that A *does* contain at least one element, and by removing the double negation we make the statement less dependent on LEM. Recall that we say that a type A is *inhabited* when we assert A itself as a proposition (i.e. we construct an element of A , usually unnamed). Thus, often when translating a classical proof into constructive logic, we replace the word “nonempty” by “inhabited” (although sometimes we must replace it instead by “merely inhabited”; see §3.7).

Similarly, it is not uncommon in classical mathematics to find unnecessary proofs by contradiction. Of course, the classical form of proof by contradiction proceeds by way of the law of double negation: we assume $\neg A$ and derive a contradiction, thereby deducing $\neg\neg A$, and thus by double negation we obtain A . However, often the derivation of a contradiction from $\neg A$ can be rephrased slightly so as to yield a direct proof of A , avoiding the need for LEM.

It is also important to note that if the goal is to prove a *negation*, then “proof by contradiction” does not involve LEM. In fact, since $\neg A$ is by definition the type $A \rightarrow 0$, by definition to prove $\neg A$ is to prove a contradiction (**0**) under the assumption of A . Similarly, the law of double negation does hold for negated propositions: $\neg\neg\neg A \rightarrow \neg A$. With practice, one learns to distinguish more carefully between negated and non-negated propositions and to notice when LEM is being used and when it is not.

Thus, contrary to how it may appear on the surface, doing mathematics “constructively” does not usually involve giving up important theorems, but rather finding the best way to state the definitions so as to make the important theorems constructively provable. That is, we may freely use the LEM when first investigating a subject, but once that subject is better understood, we can hope to refine its definitions and proofs so as to avoid that axiom. This sort of observation is even more pronounced in *homotopy* type theory, where the powerful tools of univalence and higher inductive types allow us to constructively attack many problems that traditionally would require classical reasoning. We will see several examples of this in ??.

It is also worth mentioning that even in constructive mathematics, the law of excluded middle can hold for *some* propositions. The name traditionally given to such propositions is *decidable*.

Definition 3.4.3.

- (i) A type A is called **decidable** if $A + \neg A$.

- (ii) Similarly, a type family $B : A \rightarrow \mathcal{U}$ is **decidable** if $\prod_{(a:A)}(B(a) + \neg B(a))$.
- (iii) In particular, A has **decidable equality** if $\prod_{(a,b:A)}((a = b) + \neg(a = b))$.

Thus, LEM is exactly the statement that all mere propositions are decidable, and hence so are all families of mere propositions. In particular, LEM implies that all sets (in the sense of §3.1) have decidable equality. Having decidable equality in this sense is very strong; see Theorem 7.2.5.

3.5 Subsets and propositional resizing

As another example of the usefulness of mere propositions, we discuss subsets (and more generally subtypes). Suppose $P : A \rightarrow \mathcal{U}$ is a type family, with each type $P(x)$ regarded as a proposition. Then P itself is a *predicate* on A , or a *property* of elements of A .

In set theory, whenever we have a predicate P on a set A , we may form the subset $\{x \in A \mid P(x)\}$. As mentioned briefly in §1.11, the obvious analogue in type theory is the Σ -type $\sum_{(x:A)} P(x)$. An inhabitant of $\sum_{(x:A)} P(x)$ is, of course, a pair (x, p) where $x : A$ and p is a proof of $P(x)$. However, for general P , an element $a : A$ might give rise to more than one distinct element of $\sum_{(x:A)} P(x)$, if the proposition $P(a)$ has more than one distinct proof. This is counter to the usual intuition of a subset. But if P is a *mere* proposition, then this cannot happen.

Lemma 3.5.1. *Suppose $P : A \rightarrow \mathcal{U}$ is a type family such that $P(x)$ is a mere proposition for all $x : A$. If $u, v : \sum_{(x:A)} P(x)$ are such that $\text{pr}_1(u) = \text{pr}_1(v)$, then $u = v$.*

Proof. Suppose $p : \text{pr}_1(u) = \text{pr}_1(v)$. By Theorem 2.7.2, to show $u = v$ it suffices to show $p_*(\text{pr}_2(u)) = \text{pr}_2(v)$. But $p_*(\text{pr}_2(u))$ and $\text{pr}_2(v)$ are both elements of $P(\text{pr}_1(v))$, which is a mere proposition; hence they are equal. \square

For instance, recall that in §2.4 we defined

$$(A \simeq B) := \sum_{f:A \rightarrow B} \text{isequiv}(f),$$

where each type $\text{isequiv}(f)$ was supposed to be a mere proposition. It follows that if two equivalences have equal underlying functions, then they are equal as equivalences.

Henceforth, if $P : A \rightarrow \mathcal{U}$ is a family of mere propositions (i.e. each $P(x)$ is a mere proposition), we may write

$$\{x : A \mid P(x)\} \tag{3.5.2}$$

as an alternative notation for $\sum_{(x:A)} P(x)$. (There is no technical reason not to use this notation for arbitrary P as well, but such usage could be confusing due to unintended connotations.) If A is a set, we call (3.5.2) a **subset** of A ; for general A we might call it a **subtype**. We may also refer to P itself as a *subset* or *subtype* of A ; this is actually more correct, since the type (3.5.2) in isolation doesn't remember its relationship to A .

Given such a P and $a : A$, we may write $a \in P$ or $a \in \{x : A \mid P(x)\}$ to refer to the mere proposition $P(a)$. If it holds, we may say that a is a **member** of P . Similarly, if $\{x : A \mid Q(x)\}$ is another subset of A , then we say that P is **contained** in Q , and write $P \subseteq Q$, if we have $\prod_{(x:A)}(P(x) \rightarrow Q(x))$.

As further examples of subtypes, we may define the “subuniverses” of sets and of mere propositions in a universe \mathcal{U} :

$$\begin{aligned} \text{Set}_{\mathcal{U}} &:= \{A : \mathcal{U} \mid \text{isSet}(A)\}, \\ \text{Prop}_{\mathcal{U}} &:= \{A : \mathcal{U} \mid \text{isProp}(A)\}. \end{aligned}$$

An element of $\text{Set}_{\mathcal{U}}$ is a type $A : \mathcal{U}$ together with evidence $s : \text{isSet}(A)$, and similarly for $\text{Prop}_{\mathcal{U}}$. Lemma 3.5.1 implies that $(A, s) =_{\text{Set}_{\mathcal{U}}} (B, t)$ is equivalent to $A =_{\mathcal{U}} B$ (and hence to $A \simeq B$). Thus, we will frequently abuse notation and write simply $A : \text{Set}_{\mathcal{U}}$ instead of $(A, s) : \text{Set}_{\mathcal{U}}$. We may also drop the subscript \mathcal{U} if there is no need to specify the universe in question.

Recall that for any two universes \mathcal{U}_i and \mathcal{U}_{i+1} , if $A : \mathcal{U}_i$ then also $A : \mathcal{U}_{i+1}$. Thus, for any $(A, s) : \text{Set}_{\mathcal{U}_i}$ we also have $(A, s) : \text{Set}_{\mathcal{U}_{i+1}}$, and similarly for $\text{Prop}_{\mathcal{U}_i}$, giving natural maps

$$\text{Set}_{\mathcal{U}_i} \rightarrow \text{Set}_{\mathcal{U}_{i+1}}, \quad (3.5.3)$$

$$\text{Prop}_{\mathcal{U}_i} \rightarrow \text{Prop}_{\mathcal{U}_{i+1}}. \quad (3.5.4)$$

The map (3.5.3) cannot be an equivalence, since then we could reproduce the paradoxes of self-reference that are familiar from Cantorian set theory. However, although (3.5.4) is not automatically an equivalence in the type theory we have presented so far, it is consistent to suppose that it is. That is, we may consider adding to type theory the following axiom.

Axiom 3.5.5 (Propositional resizing). *The map $\text{Prop}_{\mathcal{U}_i} \rightarrow \text{Prop}_{\mathcal{U}_{i+1}}$ is an equivalence.*

We refer to this axiom as **propositional resizing**, since it means that any mere proposition in the universe \mathcal{U}_{i+1} can be “resized” to an equivalent one in the smaller universe \mathcal{U}_i . It follows automatically if \mathcal{U}_{i+1} satisfies LEM (see Exercise 3.10). We will not assume this axiom in general, although in some places we will use it as an explicit hypothesis. It is a form of *impredicativity* for mere propositions, and by avoiding its use, the type theory is said to remain *predicative*.

In practice, what we want most frequently is a slightly different statement: that a universe \mathcal{U} under consideration contains a type which “classifies all mere propositions”. In other words, we want a type $\Omega : \mathcal{U}$ together with an Ω -indexed family of mere propositions, which contains every mere proposition up to equivalence. This statement follows from propositional resizing as stated above if \mathcal{U} is not the smallest universe \mathcal{U}_0 , since then we can define $\Omega := \text{Prop}_{\mathcal{U}_0}$.

One use for impredicativity is to define power sets. It is natural to define the **power set** of a set A to be $A \rightarrow \text{Prop}_{\mathcal{U}}$; but in the absence of impredicativity, this definition depends (even up to equivalence) on the choice of the universe \mathcal{U} . But with propositional resizing, we can define the power set to be

$$\mathcal{P}(A) := (A \rightarrow \Omega),$$

which is then independent of \mathcal{U} . See also ??.

3.6 The logic of mere propositions

We mentioned in §1.1 that in contrast to type theory, which has only one basic notion (types), set-theoretic foundations have two basic notions: sets and propositions. Thus, a classical mathematician is accustomed to manipulating these two kinds of objects separately.

It is possible to recover a similar dichotomy in type theory, with the role of the set-theoretic propositions being played by the types (and type families) that are *mere* propositions. In many cases, the logical connectives and quantifiers can be represented in this logic by simply restricting the corresponding type-former to the mere propositions. Of course, this requires knowing that the type-former in question preserves mere propositions.

Example 3.6.1. If A and B are mere propositions, so is $A \times B$. This is easy to show using the characterization of paths in products, just like Example 3.1.5 but simpler. Thus, the connective “and” preserves mere propositions.

Example 3.6.2. If A is any type and $B : A \rightarrow \mathcal{U}$ is such that for all $x : A$, the type $B(x)$ is a mere proposition, then $\prod_{(x:A)} B(x)$ is a mere proposition. The proof is just like Example 3.1.6 but simpler: given $f, g : \prod_{(x:A)} B(x)$, for any $x : A$ we have $f(x) = g(x)$ since $B(x)$ is a mere proposition. But then by function extensionality, we have $f = g$.

In particular, if B is a mere proposition, then so is $A \rightarrow B$ regardless of what A is. In even more particular, since $\mathbf{0}$ is a mere proposition, so is $\neg A \equiv (A \rightarrow \mathbf{0})$. Thus, the connectives “implies” and “not” preserve mere propositions, as does the quantifier “for all”.

On the other hand, some type formers do not preserve mere propositions. Even if A and B are mere propositions, $A + B$ will not in general be. For instance, $\mathbf{1}$ is a mere proposition, but $\mathbf{2} = \mathbf{1} + \mathbf{1}$ is not. Logically speaking, $A + B$ is a “purely constructive” sort of “or”: a witness of it contains the additional information of *which* disjunct is true. Sometimes this is very useful, but if we want a more classical sort of “or” that preserves mere propositions, we need a way to “truncate” this type into a mere proposition by forgetting this additional information.

The same issue arises with the Σ -type $\sum_{(x:A)} P(x)$. This is a purely constructive interpretation of “there exists an $x : A$ such that $P(x)$ ” which remembers the witness x , and hence is not generally a mere proposition even if each type $P(x)$ is. (Recall that we observed in §3.5 that $\sum_{(x:A)} P(x)$ can also be regarded as “the subset of those $x : A$ such that $P(x)$ ”.)

3.7 Propositional truncation

The *propositional truncation*, also called the (-1) -*truncation*, *bracket type*, or *squash type*, is an additional type former which “squashes” or “truncates” a type down to a mere proposition, forgetting all information contained in inhabitants of that type other than their existence.

More precisely, for any type A , there is a type $\|A\|$. It has two constructors:

- For any $a : A$ we have $|a| : \|A\|$.
- For any $x, y : \|A\|$, we have $x = y$.

The first constructor means that if A is inhabited, so is $\|A\|$. The second ensures that $\|A\|$ is a mere proposition; usually we leave the witness of this fact nameless.

The recursion principle of $\|A\|$ says that:

- If B is a mere proposition and we have $f : A \rightarrow B$, then there is an induced $g : \|A\| \rightarrow B$ such that $g(|a|) \equiv f(a)$ for all $a : A$.

In other words, any mere proposition which follows from (the inhabitedness of) A already follows from $\|A\|$. Thus, $\|A\|$, as a mere proposition, contains no more information than the inhabitedness of A . (There is also an induction principle for $\|A\|$, but it is not especially useful; see Exercise 3.17.)

In Exercises 3.14 and 3.15 and §6.9 we will describe some ways to construct $\|A\|$ in terms of more general things. For now, we simply assume it as an additional rule alongside those of Chapter 1.

With the propositional truncation, we can extend the “logic of mere propositions” to cover disjunction and the existential quantifier. Specifically, $\|A + B\|$ is a mere propositional version of “ A or B ”, which does not “remember” the information of which disjunct is true.

The recursion principle of truncation implies that we can still do a case analysis on $\|A + B\|$ when attempting to prove a mere proposition. That is, suppose we have an assumption $u : \|A + B\|$ and we are trying to prove a mere proposition Q . In other words, we are trying to define an

element of $\|A + B\| \rightarrow Q$. Since Q is a mere proposition, by the recursion principle for propositional truncation, it suffices to construct a function $A + B \rightarrow Q$. But now we can use case analysis on $A + B$.

Similarly, for a type family $P : A \rightarrow \mathcal{U}$, we can consider $\|\sum_{(x:A)} P(x)\|$, which is a mere propositional version of “there exists an $x : A$ such that $P(x)$ ”. As for disjunction, by combining the induction principles of truncation and Σ -types, if we have an assumption of type $\|\sum_{(x:A)} P(x)\|$, we may introduce new assumptions $x : A$ and $y : P(x)$ when attempting to prove a mere proposition. In other words, if we know that there exists some $x : A$ such that $P(x)$, but we don’t have a particular such x in hand, then we are free to make use of such an x as long as we aren’t trying to construct anything which might depend on the particular value of x . Requiring the codomain to be a mere proposition expresses this independence of the result on the witness, since all possible inhabitants of such a type must be equal.

For the purposes of set-level mathematics in Set^{Set} , where we deal mostly with sets and mere propositions, it is convenient to use the traditional logical notations to refer only to “propositionally truncated logic”.

Definition 3.7.1. We define **traditional logical notation** using truncation as follows, where P and Q denote mere propositions (or families thereof):

$$\begin{aligned}\top &:= \mathbf{1} \\ \perp &:= \mathbf{0} \\ P \wedge Q &:= P \times Q \\ P \Rightarrow Q &:= P \rightarrow Q \\ P \Leftrightarrow Q &:= P = Q \\ \neg P &:= P \rightarrow \mathbf{0} \\ P \vee Q &:= \|P + Q\| \\ \forall(x : A). P(x) &:= \prod_{x:A} P(x) \\ \exists(x : A). P(x) &:= \left\| \sum_{x:A} P(x) \right\|\end{aligned}$$

The notations \wedge and \vee are also used in homotopy theory for the smash product and the wedge of pointed spaces, which we will introduce in Chapter 6. This technically creates a potential for conflict, but no confusion will generally arise.

Similarly, when discussing subsets as in §3.5, we may use the traditional notation for intersections, unions, and complements:

$$\begin{aligned}\{x : A \mid P(x)\} \cap \{x : A \mid Q(x)\} &:= \{x : A \mid P(x) \wedge Q(x)\}, \\ \{x : A \mid P(x)\} \cup \{x : A \mid Q(x)\} &:= \{x : A \mid P(x) \vee Q(x)\}, \\ A \setminus \{x : A \mid P(x)\} &:= \{x : A \mid \neg P(x)\}.\end{aligned}$$

Of course, in the absence of LEM, the latter are not “complements” in the usual sense: we may not have $B \cup (A \setminus B) = A$ for every subset B of A .

3.8 The axiom of choice

We can now properly formulate the axiom of choice in homotopy type theory. Assume a type X and type families

$$A : X \rightarrow \mathcal{U} \quad \text{and} \quad P : \prod_{x:X} A(x) \rightarrow \mathcal{U},$$

and moreover that

- X is a set,
- $A(x)$ is a set for all $x : X$, and
- $P(x, a)$ is a mere proposition for all $x : X$ and $a : A(x)$.

The **axiom of choice** AC asserts that under these assumptions,

$$\left(\prod_{x:X} \left\| \sum_{a:A(x)} P(x, a) \right\| \right) \rightarrow \left\| \sum_{(g:\prod_{(x:X)} A(x))} \prod_{(x:X)} P(x, g(x)) \right\|. \quad (3.8.1)$$

Of course, this is a direct translation of (3.2.1) where we read “there exists $x : A$ such that $B(x)$ ” as $\left\| \sum_{(x:A)} B(x) \right\|$, so we could have written the statement in the familiar logical notation as

$$\left(\forall (x : X). \exists (a : A(x)). P(x, a) \right) \Rightarrow \left(\exists (g : \prod_{(x:X)} A(x)). \forall (x : X). P(x, g(x)) \right).$$

In particular, note that the propositional truncation appears twice. The truncation in the domain means we assume that for every x there exists some $a : A(x)$ such that $P(x, a)$, but that these values are not chosen or specified in any known way. The truncation in the codomain means we conclude that there exists some function g , but this function is not determined or specified in any known way.

In fact, because of Theorem 2.15.7, this axiom can also be expressed in a simpler form.

Lemma 3.8.2. *The axiom of choice (3.8.1) is equivalent to the statement that for any set X and any $Y : X \rightarrow \mathcal{U}$ such that each $Y(x)$ is a set, we have*

$$\left(\prod_{x:X} \left\| Y(x) \right\| \right) \rightarrow \left\| \prod_{x:X} Y(x) \right\|. \quad (3.8.3)$$

This corresponds to a well-known equivalent form of the classical axiom of choice, namely “the cartesian product of a family of nonempty sets is nonempty”.

Proof. By Theorem 2.15.7, the codomain of (3.8.1) is equivalent to

$$\left\| \prod_{(x:X)} \sum_{(a:A(x))} P(x, a) \right\|.$$

Thus, (3.8.1) is equivalent to the instance of (3.8.3) where $Y(x) := \sum_{(a:A(x))} P(x, a)$. (This is a set by Example 3.1.5 and Lemma 3.3.4.) Conversely, (3.8.3) is equivalent to the instance of (3.8.1) where $A(x) := Y(x)$ and $P(x, a) := \mathbf{1}$. Thus, the two are logically equivalent. Since both are mere propositions, by Lemma 3.3.3 they are equivalent types. \square

As with LEM, the equivalent forms (3.8.1) and (3.8.3) are not a consequence of our basic type theory, but they may consistently be assumed as axioms.

Remark 3.8.4. It is easy to show that the right side of (3.8.3) always implies the left. Since both are mere propositions, by Lemma 3.3.3 the axiom of choice is also equivalent to asking for an equivalence

$$\left(\prod_{x:X} \|Y(x)\| \right) \simeq \left\| \prod_{x:X} Y(x) \right\|$$

This illustrates a common pitfall: although dependent function types preserve mere propositions (Example 3.6.2), they do not commute with truncation: $\left\| \prod_{(x:A)} P(x) \right\|$ is not generally equivalent to $\prod_{(x:A)} \|P(x)\|$. The axiom of choice, if we assume it, says that this is true *for sets*; as we will see below, it fails in general.

The restriction in the axiom of choice to types that are sets can be relaxed to a certain extent. For instance, we may allow A and P in (3.8.1), or Y in (3.8.3), to be arbitrary type families; this results in a seemingly stronger statement that is equally consistent. We may also replace the propositional truncation by the more general n -truncations to be considered in Chapter 7, obtaining a spectrum of axioms AC_n interpolating between (3.8.1), which we call simply AC (or AC_{-1} for emphasis), and Theorem 2.15.7, which we shall call AC_∞ . See also Exercises 7.8 and 7.10. However, observe that we cannot relax the requirement that X be a set.

Lemma 3.8.5. *There exists a type X and a family $Y : X \rightarrow \mathcal{U}$ such that each $Y(x)$ is a set, but such that (3.8.3) is false.*

Proof. Define $X := \sum_{(A:\mathcal{U})} \|\mathbf{2} = A\|$, and let $x_0 := (\mathbf{2}, |\text{refl}_2|) : X$. Then by the identification of paths in Σ -types, the fact that $\|\mathbf{2} = \mathbf{2}\|$ is a mere proposition, and univalence, for any $(A, p), (B, q) : X$ we have $((A, p) =_X (B, q)) \simeq (A \simeq B)$. In particular, $(x_0 =_X x_0) \simeq (\mathbf{2} \simeq \mathbf{2})$, so as in Example 3.1.9, X is not a set.

On the other hand, if $(A, p) : X$, then A is a set; this follows by induction on truncation for $p : \|\mathbf{2} = A\|$ and the fact that $\mathbf{2}$ is a set. Since $A \simeq B$ is a set whenever A and B are, it follows that $x_1 =_X x_2$ is a set for any $x_1, x_2 : X$, i.e. X is a 1-type. In particular, if we define $Y : X \rightarrow \mathcal{U}$ by $Y(x) := (x_0 = x)$, then each $Y(x)$ is a set.

Now by definition, for any $(A, p) : X$ we have $\|\mathbf{2} = A\|$, and hence $\|x_0 = (A, p)\|$. Thus, we have $\prod_{(x:X)} \|Y(x)\|$. If (3.8.3) held for this X and Y , then we would also have $\left\| \prod_{(x:X)} Y(x) \right\|$. Since we are trying to derive a contradiction (0), which is a mere proposition, we may assume $\prod_{(x:X)} Y(x)$, i.e. that $\prod_{(x:X)} (x_0 = x)$. But this implies X is a mere proposition, and hence a set, which is a contradiction. \square

3.9 The principle of unique choice

The following observation is trivial, but very useful.

Lemma 3.9.1. *If P is a mere proposition, then $P \simeq \|P\|$.*

Proof. Of course, we have $P \rightarrow \|P\|$ by definition. And since P is a mere proposition, the universal property of $\|P\|$ applied to $\text{id}_P : P \rightarrow P$ yields $\|P\| \rightarrow P$. These functions are quasi-inverses by Lemma 3.3.3. \square

Among its important consequences is the following.

Corollary 3.9.2 (The principle of unique choice). *Suppose a type family $P : A \rightarrow \mathcal{U}$ such that*

- (i) For each x , the type $P(x)$ is a mere proposition, and
- (ii) For each x we have $\|P(x)\|$.

Then we have $\prod_{(x:A)} P(x)$.

Proof. Immediate from the two assumptions and the previous lemma. \square

The corollary also encapsulates a very useful technique of reasoning. Namely, suppose we know that $\|A\|$, and we want to use this to construct an element of some other type B . We would like to use an element of A in our construction of an element of B , but this is allowed only if B is a mere proposition, so that we can apply the induction principle for the propositional truncation $\|B\|$; the most we could hope to do in general is to show $\|B\|$. Instead, we can extend B with additional data which characterizes *uniquely* the object we wish to construct. Specifically, we define a predicate $Q : B \rightarrow \mathcal{U}$ such that $\sum_{(x:B)} Q(x)$ is a mere proposition. Then from an element of A we construct an element $b : B$ such that $Q(b)$, hence from $\|A\|$ we can construct $\|\sum_{(x:B)} Q(x)\|$, and because $\|\sum_{(x:B)} Q(x)\|$ is equivalent to $\sum_{(x:B)} Q(x)$ an element of B may be projected from it. An example can be found in Exercise 3.19.

A similar issue arises in set-theoretic mathematics, although it manifests slightly differently. If we are trying to define a function $f : A \rightarrow B$, and depending on an element $a : A$ we are able to prove mere existence of some $b : B$, we are not done yet because we need to actually pinpoint an element of B , not just prove its existence. One option is of course to refine the argument to unique existence of $b : B$, as we did in type theory. But in set theory the problem can often be avoided more simply by an application of the axiom of choice, which picks the required elements for us. In homotopy type theory, however, quite apart from any desire to avoid choice, the available forms of choice are simply less applicable, since they require that the domain of choice be a *set*. Thus, if A is not a set (such as perhaps a universe \mathcal{U}), there is no consistent form of choice that will allow us to simply pick an element of B for each $a : A$ to use in defining $f(a)$.

3.10 When are propositions truncated?

At first glance, it may seem that the truncated versions of $+$ and Σ are actually closer to the informal mathematical meaning of “or” and “there exists” than the untruncated ones. Certainly, they are closer to the *precise* meaning of “or” and “there exists” in the first-order logic which underlies formal set theory, since the latter makes no attempt to remember any witnesses to the truth of propositions. However, it may come as a surprise to realize that the practice of *informal* mathematics is often more accurately described by the untruncated forms.

For example, consider a statement like “every prime number is either 2 or odd”. The working mathematician feels no compunction about using this fact not only to prove *theorems* about prime numbers, but also to perform *constructions* on prime numbers, perhaps doing one thing in the case of 2 and another in the case of an odd prime. The end result of the construction is not merely the truth of some statement, but a piece of data which may depend on the parity of the prime number. Thus, from a type-theoretic perspective, such a construction is naturally phrased using the induction principle for the coproduct type “ $(p = 2) + (p \text{ is odd})$ ”, not its propositional truncation.

Admittedly, this is not an ideal example, since “ $p = 2$ ” and “ p is odd” are mutually exclusive, so that $(p = 2) + (p \text{ is odd})$ is in fact already a mere proposition and hence equivalent to its truncation (see Exercise 3.7). More compelling examples come from the existential quantifier. It

is not uncommon to prove a theorem of the form “there exists an x such that …” and then refer later on to “the x constructed in Theorem Y” (note the definite article). Moreover, when deriving further properties of this x , one may use phrases such as “by the construction of x in the proof of Theorem Y”.

A very common example is “ A is isomorphic to B ”, which strictly speaking means only that there exists *some* isomorphism between A and B . But almost invariably, when proving such a statement, one exhibits a specific isomorphism or proves that some previously known map is an isomorphism, and it often matters later on what particular isomorphism was given.

Set-theoretically trained mathematicians often feel a twinge of guilt at such “abuses of language”. We may attempt to apologize for them, expunge them from final drafts, or weasel out of them with vague words like “canonical”. The problem is exacerbated by the fact that in formalized set theory, there is technically no way to “construct” objects at all — we can only prove that an object with certain properties exists. Untruncated logic in type theory thus captures some common practices of informal mathematics that the set theoretic reconstruction obscures. (This is similar to how the univalence axiom validates the common, but formally unjustified, practice of identifying isomorphic objects.)

On the other hand, sometimes truncated logic is essential. We have seen this in the statements of LEM and AC; some other examples will appear later on in the book. Thus, we are faced with the problem: when writing informal type theory, what should we mean by the words “or” and “there exists” (along with common synonyms such as “there is” and “we have”)?

A universal consensus may not be possible. Perhaps depending on the sort of mathematics being done, one convention or the other may be more useful — or, perhaps, the choice of convention may be irrelevant. In this case, a remark at the beginning of a mathematical paper may suffice to inform the reader of the linguistic conventions in use therein. However, even after one overall convention is chosen, the other sort of logic will usually arise at least occasionally, so we need a way to refer to it. More generally, one may consider replacing the propositional truncation with another operation on types that behaves similarly, such as the double negation operation $A \mapsto \neg\neg A$, or the n -truncations to be considered in Chapter 7. As an experiment in exposition, in what follows we will occasionally use *adverbs* to denote the application of such “modalities” as propositional truncation.

For instance, if untruncated logic is the default convention, we may use the adverb **merely** to denote propositional truncation. Thus the phrase

“there merely exists an $x : A$ such that $P(x)$ ”

indicates the type $\|\sum_{(x:A)} P(x)\|$. Similarly, we will say that a type A is **merely inhabited** to mean that its propositional truncation $\|A\|$ is inhabited (i.e. that we have an unnamed element of it). Note that this is a *definition* of the adverb “merely” as it is to be used in our informal mathematical English, in the same way that we define nouns like “group” and “ring”, and adjectives like “regular” and “normal”, to have precise mathematical meanings. We are not claiming that the dictionary definition of “merely” refers to propositional truncation; the choice of word is meant only to remind the mathematician reader that a mere proposition contains “merely” the information of a truth value and nothing more.

On the other hand, if truncated logic is the current default convention, we may use an adverb such as **purely** or **constructively** to indicate its absence, so that

“there purely exists an $x : A$ such that $P(x)$ ”

would denote the type $\sum_{(x:A)} P(x)$. We may also use “purely” or “actually” just to emphasize the absence of truncation, even when that is the default convention.

In this book we will continue using untruncated logic as the default convention, for a number of reasons.

- (1) We want to encourage the newcomer to experiment with it, rather than sticking to truncated logic simply because it is more familiar.
- (2) Using truncated logic as the default in type theory suffers from the same sort of “abuse of language” problems as set-theoretic foundations, which untruncated logic avoids. For instance, our definition of “ $A \simeq B$ ” as the type of equivalences between A and B , rather than its propositional truncation, means that to prove a theorem of the form “ $A \simeq B$ ” is literally to construct a particular such equivalence. This specific equivalence can then be referred to later on.
- (3) We want to emphasize that the notion of “mere proposition” is not a fundamental part of type theory. As we will see in Chapter 7, mere propositions are just the second rung on an infinite ladder, and there are also many other modalities not lying on this ladder at all.
- (4) Many statements that classically are mere propositions are no longer so in homotopy type theory. Of course, foremost among these is equality.
- (5) On the other hand, one of the most interesting observations of homotopy type theory is that a surprising number of types are *automatically* mere propositions, or can be slightly modified to become so, without the need for any truncation. (See Lemma 3.3.5, Chapters 4 and 7, and ????.) Thus, although these types contain no data beyond a truth value, we can nevertheless use them to construct untruncated objects, since there is no need to use the induction principle of propositional truncation. This useful fact is more clumsy to express if propositional truncation is applied to all statements by default.
- (6) Finally, truncations are not very useful for most of the mathematics we will be doing in this book, so it is simpler to notate them explicitly when they occur.

3.11 Contractibility

In Lemma 3.3.2 we observed that a mere proposition which is inhabited must be equivalent to 1 , and it is not hard to see that the converse also holds. A type with this property is called *contractible*. Another equivalent definition of contractibility, which is also sometimes convenient, is the following.

Definition 3.11.1. A type A is **contractible**, or a **singleton**, if there is $a : A$, called the **center of contraction**, such that $a = x$ for all $x : A$. We denote the specified path $a = x$ by contr_x .

In other words, the type $\text{isContr}(A)$ is defined to be

$$\text{isContr}(A) := \sum_{(a:A)} \prod_{(x:A)} (a = x).$$

Note that under the usual propositions-as-types reading, we can pronounce $\text{isContr}(A)$ as “ A contains exactly one element”, or more precisely “ A contains an element, and every element of A is equal to that element”.

Remark 3.11.2. We can also pronounce $\text{isContr}(A)$ more topologically as “there is a point $a : A$ such that for all $x : A$ there exists a path from a to x ”. Note that to a classical ear, this sounds like a definition of *connectedness* rather than contractibility. The point is that the meaning of “there exists” in this sentence is a continuous/natural one.

A better way to express connectedness would be $\sum_{(a:A)} \prod_{(x:A)} \|a = x\|$. This is indeed correct if A is assumed to be pointed — see the remark after Lemma 7.5.11 — but in general a type can be connected without being pointed. In §7.5 we will define connectedness as the $n = 0$ case of a general notion of n -connectedness, and in Exercise 7.6 the reader is asked to show that this definition is equivalent to having both $\|A\|$ and $\prod_{(x,y:A)} \|x = y\|$.

Lemma 3.11.3. *For a type A , the following are logically equivalent.*

- (i) *A is contractible in the sense of Definition 3.11.1.*
- (ii) *A is a mere proposition, and there is a point $a : A$.*
- (iii) *A is equivalent to $\mathbf{1}$.*

Proof. If A is contractible, then it certainly has a point $a : A$ (the center of contraction), while for any $x, y : A$ we have $x = a = y$; thus A is a mere proposition. Conversely, if we have $a : A$ and A is a mere proposition, then for any $x : A$ we have $x = a$; thus A is contractible. And we showed (ii) \Rightarrow (iii) in Lemma 3.3.2, while the converse follows since $\mathbf{1}$ easily has property (ii). \square

Lemma 3.11.4. *For any type A , the type $\text{isContr}(A)$ is a mere proposition.*

Proof. Suppose given $c, c' : \text{isContr}(A)$. We may assume $c \equiv (a, p)$ and $c' \equiv (a', p')$ for $a, a' : A$ and $p : \prod_{(x:A)} (a = x)$ and $p' : \prod_{(x:A)} (a' = x)$. By the characterization of paths in Σ -types, to show $c = c'$ it suffices to exhibit $q : a = a'$ such that $q_*(p) = p'$. We choose $q := p(a')$. Now since A is contractible (by c or c'), by Lemma 3.11.3 it is a mere proposition. Hence, by Lemma 3.3.4 and Example 3.6.2, so is $\prod_{(x:A)} (a' = x)$; thus $q_*(p) = p'$ is automatic. \square

Corollary 3.11.5. *If A is contractible, then so is $\text{isContr}(A)$.*

Proof. By Lemma 3.11.4 and Lemma 3.11.3(ii). \square

Like mere propositions, contractible types are preserved by many type constructors. For instance, we have:

Lemma 3.11.6. *If $P : A \rightarrow \mathcal{U}$ is a type family such that each $P(a)$ is contractible, then $\prod_{(x:A)} P(x)$ is contractible.*

Proof. By Example 3.6.2, $\prod_{(x:A)} P(x)$ is a mere proposition since each $P(x)$ is. But it also has an element, namely the function sending each $x : A$ to the center of contraction of $P(x)$. Thus by Lemma 3.11.3(ii), $\prod_{(x:A)} P(x)$ is contractible. \square

(In fact, the statement of Lemma 3.11.6 is equivalent to the function extensionality axiom. See §4.9.)

Of course, if A is equivalent to B and A is contractible, then so is B . More generally, it suffices for B to be a *retract* of A . By definition, a **retraction** is a function $r : A \rightarrow B$ such that there exists a function $s : B \rightarrow A$, called its **section**, and a homotopy $\epsilon : \prod_{(y:B)} (r(s(y)) = y)$; then we say that B is a **retract** of A .

Lemma 3.11.7. *If B is a retract of A , and A is contractible, then so is B .*

Proof. Let $a_0 : A$ be the center of contraction. We claim that $b_0 := r(a_0) : B$ is a center of contraction for B . Let $b : B$; we need a path $b = b_0$. But we have $\epsilon_b : r(s(b)) = b$ and $\text{contr}_{s(b)} : s(b) = a_0$, so by composition

$$\epsilon_b^{-1} \bullet r(\text{contr}_{s(b)}) : b = r(a_0) \equiv b_0. \quad \square$$

Contractible types may not seem very interesting, since they are all equivalent to **1**. One reason the notion is useful is that sometimes a collection of individually nontrivial data will collectively form a contractible type. An important example is the space of paths with one free endpoint. As we will see in §5.8, this fact essentially encapsulates the based path induction principle for identity types.

Lemma 3.11.8. *For any A and any $a : A$, the type $\sum_{(x:A)}(a = x)$ is contractible.*

Proof. We choose as center the point (a, refl_a) . Now suppose $(x, p) : \sum_{(x:A)}(a = x)$; we must show $(a, \text{refl}_a) = (x, p)$. By the characterization of paths in Σ -types, it suffices to exhibit $q : a = x$ such that $q_*(\text{refl}_a) = p$. But we can take $q := p$, in which case $q_*(\text{refl}_a) = p$ follows from the characterization of transport in path types. \square

When this happens, it can allow us to simplify a complicated construction up to equivalence, using the informal principle that contractible data can be freely ignored. This principle consists of many lemmas, most of which we leave to the reader; the following is an example.

Lemma 3.11.9. *Let $P : A \rightarrow \mathcal{U}$ be a type family.*

- (i) *If each $P(x)$ is contractible, then $\sum_{(x:A)} P(x)$ is equivalent to A .*
- (ii) *If A is contractible with center a , then $\sum_{(x:A)} P(x)$ is equivalent to $P(a)$.*

Proof. In the situation of (i), we show that $\text{pr}_1 : \sum_{(x:A)} P(x) \rightarrow A$ is an equivalence. For quasi-inverse we define $g(x) := (x, c_x)$ where c_x is the center of $P(x)$. The composite $\text{pr}_1 \circ g$ is obviously id_A , whereas the opposite composite is homotopic to the identity by using the contractions of each $P(x)$.

We leave the proof of (ii) to the reader (see Exercise 3.20). \square

Another reason contractible types are interesting is that they extend the ladder of n -types mentioned in §3.1 downwards one more step.

Lemma 3.11.10. *A type A is a mere proposition if and only if for all $x, y : A$, the type $x =_A y$ is contractible.*

Proof. For “if”, we simply observe that any contractible type is inhabited. For “only if”, we observed in §3.3 that every mere proposition is a set, so that each type $x =_A y$ is a mere proposition. But it is also inhabited (since A is a mere proposition), and hence by Lemma 3.11.3(ii) it is contractible. \square

Thus, contractible types may also be called **(−2)-types**. They are the bottom rung of the ladder of n -types, and will be the base case of the recursive definition of n -types in Chapter 7.

Notes

The fact that it is possible to define sets, mere propositions, and contractible types in type theory, with all higher homotopies automatically taken care of as in §§3.1, 3.3 and 3.11, was first observed by Voevodsky. In fact, he defined the entire hierarchy of n -types by induction, as we will do in Chapter 7.

Theorem 3.2.2 and Corollary 3.2.7 rely in essence on a classical theorem of Hedberg, which we will prove in §7.2. The implication that the propositions-as-types form of LEM contradicts univalence was observed by Martín Escardó on the AGDA mailing list. The proof we have given of Theorem 3.2.2 is due to Thierry Coquand.

The propositional truncation was introduced in the extensional type theory of NUPRL in 1983 by Constable [Con85] as an application of “subset” and “quotient” types. What is here called the “propositional truncation” was called “squashing” in the NUPRL type theory [CAB⁺86]. Rules characterizing the propositional truncation directly, still in extensional type theory, were given in [AB04]. The intensional version in homotopy type theory was constructed by Voevodsky using an impredicative quantification, and later by Lumsdaine using higher inductive types (see §6.9).

Voevodsky [Voe12] has proposed resizing rules of the kind considered in §3.5. These are clearly related to the notorious *axiom of reducibility* proposed by Russell in his and Whitehead’s *Principia Mathematica* [WR27].

The adverb “purely” as used to refer to untruncated logic is a reference to the use of monadic modalities to model effects in programming languages; see §7.7 and the Notes to Chapter 7.

There are many different ways in which logic can be treated relative to type theory. For instance, in addition to the plain propositions-as-types logic described in §1.11, and the alternative which uses mere propositions only as described in §3.6, one may introduce a separate “sort” of propositions, which behave somewhat like types but are not identified with them. This is the approach taken in logic enriched type theory [AG02] and in some presentations of the internal languages of toposes and related categories (e.g. [Jac99, Joh02]), as well as in the proof assistant COQ. Such an approach is more general, but less powerful. For instance, the principle of unique choice (§3.9) fails in the category of so-called setoids in COQ [Spi11], in logic enriched type theory [AG02], and in minimal type theory [MS05]. Thus, the univalence axiom makes our type theory behave more like the internal logic of a topos; see also ??.

Martin-Löf [ML06] provides a discussion on the history of axioms of choice. Of course, constructive and intuitionistic mathematics has a long and complicated history, which we will not delve into here; see for instance [TvD88a, TvD88b].

Exercises

Exercise 3.1. Prove that if $A \simeq B$ and A is a set, then so is B .

Exercise 3.2. Prove that if A and B are sets, then so is $A + B$.

Exercise 3.3. Prove that if A is a set and $B : A \rightarrow \mathcal{U}$ is a type family such that $B(x)$ is a set for all $x : A$, then $\sum_{(x:A)} B(x)$ is a set.

Exercise 3.4. Show that A is a mere proposition if and only if $A \rightarrow A$ is contractible.

Exercise 3.5. Show that $\text{isProp}(A) \simeq (A \rightarrow \text{isContr}(A))$.

Exercise 3.6. Show that if A is a mere proposition, then so is $A + (\neg A)$. Thus, there is no need to insert a propositional truncation in (3.4.1).

Exercise 3.7. More generally, show that if A and B are mere propositions and $\neg(A \times B)$, then $A + B$ is also a mere proposition.

Exercise 3.8. Assuming that some type $\text{isequiv}(f)$ satisfies conditions (i)–(iii) of §2.4, show that the type $\|\text{qinv}(f)\|$ satisfies the same conditions and is equivalent to $\text{isequiv}(f)$.

Exercise 3.9. Show that if LEM holds, then the type $\text{Prop} := \sum_{(A:\mathcal{U})} \text{isProp}(A)$ is equivalent to $\mathbf{2}$.

Exercise 3.10. Show that if \mathcal{U}_{i+1} satisfies LEM, then the canonical inclusion $\text{Prop}_{\mathcal{U}_i} \rightarrow \text{Prop}_{\mathcal{U}_{i+1}}$ is an equivalence.

Exercise 3.11. Show that it is not the case that for all $A : \mathcal{U}$ we have $\|A\| \rightarrow A$. (However, there can be particular types for which $\|A\| \rightarrow A$. Exercise 3.8 implies that $\text{qinv}(f)$ is such.)

Exercise 3.12. Show that if LEM holds, then for all $A : \mathcal{U}$ we have $\|(\|A\| \rightarrow A)\|$. (This property is a very simple form of the axiom of choice, which can fail in the absence of LEM; see [KECA13].)

Exercise 3.13. We showed in Corollary 3.2.7 that the following naive form of LEM is inconsistent with univalence:

$$\prod_{A:\mathcal{U}} (A + (\neg A))$$

In the absence of univalence, this axiom is consistent. However, show that it implies the axiom of choice (3.8.1).

Exercise 3.14. Show that assuming LEM, the double negation $\neg\neg A$ has the same recursion principle as the propositional truncation $\|A\|$ but with a propositional computation rule rather than a judgmental one. In other words, prove that assuming LEM, if B is a mere proposition and we have $f : A \rightarrow B$, then there is an induced $g : \neg\neg A \rightarrow B$ such that $g(|a|) = f(a)$ for all $a : A$. Deduce that (assuming LEM) we have $\neg\neg A \simeq \|A\|$. Thus, under LEM, the propositional truncation can be defined rather than taken as a separate type former.

Exercise 3.15. Show that if we assume propositional resizing as in §3.5, then the type

$$\prod_{P:\text{Prop}} ((A \rightarrow P) \rightarrow P)$$

has the same recursion principle as $\|A\|$, with the same judgmental computation rule. Thus, we can also define the propositional truncation in this case.

Exercise 3.16. Assuming LEM, show that double negation commutes with universal quantification of mere propositions over sets. That is, show that if X is a set and each $Y(x)$ is a mere proposition, then LEM implies

$$\left(\prod_{x:X} \neg\neg Y(x) \right) \simeq \left(\neg\neg \prod_{x:X} Y(x) \right). \quad (3.11.11)$$

Observe that if we assume instead that each $Y(x)$ is a set, then (3.11.11) becomes equivalent to the axiom of choice (3.8.3).

Exercise 3.17. Show that the rules for the propositional truncation given in §3.7 are sufficient to imply the following induction principle: for any type family $B : \|A\| \rightarrow \mathcal{U}$ such that each $B(x)$ is a mere proposition, if for every $a : A$ we have $B(|a|)$, then for every $x : \|A\|$ we have $B(x)$.

Exercise 3.18. Show that the law of excluded middle (3.4.1) and the law of double negation (3.4.2) are logically equivalent.

Exercise 3.19. Suppose $P : \mathbb{N} \rightarrow \mathcal{U}$ is a decidable family (see Definition 3.4.3(ii)) of mere propositions. Prove that

$$\left\| \sum_{n:\mathbb{N}} P(n) \right\| \rightarrow \sum_{n:\mathbb{N}} P(n).$$

Exercise 3.20. Prove Lemma 3.11.9(ii): if A is contractible with center a , then $\sum_{(x:A)} P(x)$ is equivalent to $P(a)$.

Exercise 3.21. Prove that $\text{isProp}(P) \simeq (P \simeq \|P\|)$.

Exercise 3.22. As in classical set theory, the finite version of the axiom of choice is a theorem. Prove that the axiom of choice (3.8.1) holds when X is a finite type $\text{Fin}(n)$ (as defined in Exercise 1.9).

Exercise 3.23. Show that the conclusion of Exercise 3.19 is true if $P : \mathbb{N} \rightarrow \mathcal{U}$ is any decidable family.

Chapter 4

Equivalences

We now study in more detail the notion of *equivalence of types* that was introduced briefly in §2.4. Specifically, we will give several different ways to define a type $\text{isequiv}(f)$ having the properties mentioned there. Recall that we wanted $\text{isequiv}(f)$ to have the following properties, which we restate here:

- (i) $\text{qinv}(f) \rightarrow \text{isequiv}(f)$.
- (ii) $\text{isequiv}(f) \rightarrow \text{qinv}(f)$.
- (iii) $\text{isequiv}(f)$ is a mere proposition.

Here $\text{qinv}(f)$ denotes the type of quasi-inverses to f :

$$\sum_{g:B \rightarrow A} ((f \circ g \sim \text{id}_B) \times (g \circ f \sim \text{id}_A)).$$

By function extensionality, it follows that $\text{qinv}(f)$ is equivalent to the type

$$\sum_{g:B \rightarrow A} ((f \circ g = \text{id}_B) \times (g \circ f = \text{id}_A)).$$

We will define three different types having properties (i)–(iii), which we call

- half adjoint equivalences,
- bi-invertible maps, and
- contractible functions.

We will also show that all these types are equivalent. These names are intentionally somewhat cumbersome, because after we know that they are all equivalent and have properties (i)–(iii), we will revert to saying simply “equivalence” without needing to specify which particular definition we choose. But for purposes of the comparisons in this chapter, we need different names for each definition.

Before we examine the different notions of equivalence, however, we give a little more explanation of why a different concept than quasi-invertibility is needed.

4.1 Quasi-inverses

We have said that $\text{qinv}(f)$ is unsatisfactory because it is not a mere proposition, whereas we would rather that a given function could “be an equivalence” in at most one way. However, we have given no evidence that $\text{qinv}(f)$ is not a mere proposition. In this section we exhibit a specific counterexample.

Lemma 4.1.1. *If $f : A \rightarrow B$ is such that $\text{qinv}(f)$ is inhabited, then*

$$\text{qinv}(f) \simeq \left(\prod_{x:A} (x = x) \right).$$

Proof. By assumption, f is an equivalence; that is, we have $e : \text{isequiv}(f)$ and so $(f, e) : A \simeq B$. By univalence, $\text{idtoeqv} : (A = B) \rightarrow (A \simeq B)$ is an equivalence, so we may assume that (f, e) is of the form $\text{idtoeqv}(p)$ for some $p : A = B$. Then by path induction, we may assume p is refl_A , in which case f is id_A . Thus we are reduced to proving $\text{qinv}(\text{id}_A) \simeq (\prod_{(x:A)} (x = x))$. Now by definition we have

$$\text{qinv}(\text{id}_A) \equiv \sum_{g:A \rightarrow A} ((g \sim \text{id}_A) \times (g \sim \text{id}_A)).$$

By function extensionality, this is equivalent to

$$\sum_{g:A \rightarrow A} ((g = \text{id}_A) \times (g = \text{id}_A)).$$

And by Exercise 2.10, this is equivalent to

$$\sum_{h:\sum_{(g:A \rightarrow A)} (g = \text{id}_A)} (\text{pr}_1(h) = \text{id}_A)$$

However, by Lemma 3.11.8, $\sum_{(g:A \rightarrow A)} (g = \text{id}_A)$ is contractible with center $(\text{id}_A, \text{refl}_{\text{id}_A})$; therefore by Lemma 3.11.9 this type is equivalent to $\text{id}_A = \text{id}_A$. And by function extensionality, $\text{id}_A = \text{id}_A$ is equivalent to $\prod_{(x:A)} x = x$. \square

We remark that Exercise 4.3 asks for a proof of the above lemma which avoids univalence.

Thus, what we need is some A which admits a nontrivial element of $\prod_{(x:A)} (x = x)$. Thinking of A as a higher groupoid, an inhabitant of $\prod_{(x:A)} (x = x)$ is a natural transformation from the identity functor of A to itself. Such transformations are said to form the **center of a category**, since the naturality axiom requires that they commute with all morphisms. Classically, if A is simply a group regarded as a one-object groupoid, then this yields precisely its center in the usual group-theoretic sense. This provides some motivation for the following.

Lemma 4.1.2. *Suppose we have a type A with $a : A$ and $q : a = a$ such that*

- (i) *The type $a = a$ is a set.*
- (ii) *For all $x : A$ we have $\|a = x\|$.*
- (iii) *For all $p : a = a$ we have $p \cdot q = q \cdot p$.*

Then there exists $f : \prod_{(x:A)} (x = x)$ with $f(a) = q$.

Proof. Let $g : \prod_{(x:A)} \|a = x\|$ be as given by (ii). First we observe that each type $x =_A y$ is a set. For since being a set is a mere proposition, we may apply the induction principle of propositional truncation, and assume that $g(x) = |p|$ and $g(y) = |p'|$ for $p : a = x$ and $p' : a = y$. In this case, composing with p and p'^{-1} yields an equivalence $(x = y) \simeq (a = a)$. But $(a = a)$ is a set by (i), so $(x = y)$ is also a set.

Now, we would like to define f by assigning to each x the path $g(x)^{-1} \cdot q \cdot g(x)$, but this does not work because $g(x)$ does not inhabit $a = x$ but rather $\|a = x\|$, and the type $(x = x)$ may not be a mere proposition, so we cannot use induction on propositional truncation. Instead we can apply the technique mentioned in §3.9: we characterize uniquely the object we wish to construct. Let us define, for each $x : A$, the type

$$B(x) := \sum_{(r:x=x)} \prod_{(s:a=x)} (r = s^{-1} \cdot q \cdot s).$$

We claim that $B(x)$ is a mere proposition for each $x : A$. Since this claim is itself a mere proposition, we may again apply induction on truncation and assume that $g(x) = |p|$ for some $p : a = x$. Now suppose given (r, h) and (r', h') in $B(x)$; then we have

$$h(p) \cdot h'(p)^{-1} : r = r'.$$

It remains to show that h is identified with h' when transported along this equality, which by transport in identity types and function types (§§2.9 and 2.11), reduces to showing

$$h(s) = h(p) \cdot h'(p)^{-1} \cdot h'(s)$$

for any $s : a = x$. But each side of this is an equality between elements of $(x = x)$, so it follows from our above observation that $(x = x)$ is a set.

Thus, each $B(x)$ is a mere proposition; we claim that $\prod_{(x:A)} B(x)$. Given $x : A$, we may now invoke the induction principle of propositional truncation to assume that $g(x) = |p|$ for $p : a = x$. We define $r := p^{-1} \cdot q \cdot p$; to inhabit $B(x)$ it remains to show that for any $s : a = x$ we have $r = s^{-1} \cdot q \cdot s$. Manipulating paths, this reduces to showing that $q \cdot (p \cdot s^{-1}) = (p \cdot s^{-1}) \cdot q$. But this is just an instance of (iii). \square

Theorem 4.1.3. *There exist types A and B and a function $f : A \rightarrow B$ such that $\text{qinv}(f)$ is not a mere proposition.*

Proof. It suffices to exhibit a type X such that $\prod_{(x:X)} (x = x)$ is not a mere proposition. Define $X := \sum_{(A:\mathcal{U})} \|\mathbf{2} = A\|$, as in the proof of Lemma 3.8.5. It will suffice to exhibit an $f : \prod_{(x:X)} (x = x)$ which is unequal to $\lambda x. \text{refl}_x$.

Let $a := (\mathbf{2}, |\text{refl}_2|) : X$, and let $q : a = a$ be the path corresponding to the nonidentity equivalence $e : \mathbf{2} \simeq \mathbf{2}$ defined by $e(0_2) := 1_2$ and $e(1_2) := 0_2$. We would like to apply Lemma 4.1.2 to build an f . By definition of X , equalities in subset types (§3.5), and univalence, we have $(a = a) \simeq (\mathbf{2} \simeq \mathbf{2})$, which is a set, so (i) holds. Similarly, by definition of X and equalities in subset types we have (ii). Finally, Exercise 2.13 implies that every equivalence $\mathbf{2} \simeq \mathbf{2}$ is equal to either id_2 or e , so we can show (iii) by a four-way case analysis.

Thus, we have $f : \prod_{(x:X)} (x = x)$ such that $f(a) = q$. Since e is not equal to id_2 , q is not equal to refl_a , and thus f is not equal to $\lambda x. \text{refl}_x$. Therefore, $\prod_{(x:X)} (x = x)$ is not a mere proposition. \square

More generally, Lemma 4.1.2 implies that any “Eilenberg–Mac Lane space” $K(G, 1)$, where G is a nontrivial abelian group, will provide a counterexample; see ???. The type X we used turns

out to be equivalent to $K(\mathbb{Z}_2, 1)$. In Chapter 6 we will see that the circle $S^1 = K(\mathbb{Z}, 1)$ is another easy-to-describe example.

We now move on to describing better notions of equivalence.

4.2 Half adjoint equivalences

In §4.1 we concluded that $\text{qinv}(f)$ is equivalent to $\prod_{(x:A)}(x = x)$ by discarding a contractible type. Roughly, the type $\text{qinv}(f)$ contains three data g , η , and ϵ , of which two (g and η) could together be seen to be contractible when f is an equivalence. The problem is that removing these data left one remaining (ϵ). In order to solve this problem, the idea is to add one *additional* datum which, together with ϵ , forms a contractible type.

Definition 4.2.1. A function $f : A \rightarrow B$ is a **half adjoint equivalence** if there are $g : B \rightarrow A$ and homotopies $\eta : g \circ f \sim \text{id}_A$ and $\epsilon : f \circ g \sim \text{id}_B$ such that there exists a homotopy

$$\tau : \prod_{x:A} f(\eta x) = \epsilon(fx).$$

Thus we have a type $\text{ishae}(f)$, defined to be

$$\sum_{(g:B \rightarrow A)} \sum_{(\eta:g \circ f \sim \text{id}_A)} \sum_{(\epsilon:f \circ g \sim \text{id}_B)} \prod_{(x:A)} f(\eta x) = \epsilon(fx).$$

Note that in the above definition, the coherence condition relating η and ϵ only involves f . We might consider instead an analogous coherence condition involving g :

$$\nu : \prod_{y:B} g(\epsilon y) = \eta(gy)$$

and a resulting analogous definition $\text{ishae}'(f)$.

Fortunately, it turns out each of the conditions implies the other one:

Lemma 4.2.2. For functions $f : A \rightarrow B$ and $g : B \rightarrow A$ and homotopies $\eta : g \circ f \sim \text{id}_A$ and $\epsilon : f \circ g \sim \text{id}_B$, the following conditions are logically equivalent:

- $\prod_{(x:A)} f(\eta x) = \epsilon(fx)$
- $\prod_{(y:B)} g(\epsilon y) = \eta(gy)$

Proof. It suffices to show one direction; the other one is obtained by replacing A , f , and η by B , g , and ϵ respectively. Let $\tau : \prod_{(x:A)} f(\eta x) = \epsilon(fx)$. Fix $y : B$. Using naturality of ϵ and applying g , we get the following commuting diagram of paths:

$$\begin{array}{ccc} gfgfgy & \xrightarrow{\quad gfg(\epsilon y) \quad} & gfgy \\ \parallel & & \parallel \\ gfgy & \xrightarrow[\quad g(\epsilon y) \quad]{} & gy \end{array}$$

Using $\tau(gy)$ on the left side of the diagram gives us

$$\begin{array}{ccc} gfgfgy & \xrightarrow{\quad gfg(\epsilon y) \quad} & gfgy \\ \parallel & & \parallel \\ gf(\eta(gy)) & \xrightarrow[\quad g(\epsilon y) \quad]{} & gy \end{array}$$

Using the commutativity of η with $g \circ f$ (Corollary 2.4.4), we have

$$\begin{array}{ccc} gfgfgy & \xrightarrow{gfg(\epsilon y)} & gfgy \\ \eta(gfgy) \parallel & & \parallel g(\epsilon y) \\ gfgy & \xrightarrow{g(\epsilon y)} & gy \end{array}$$

However, by naturality of η we also have

$$\begin{array}{ccc} gfgfgy & \xrightarrow{gfg(\epsilon y)} & gfgy \\ \eta(gfgy) \parallel & & \parallel \eta(gy) \\ gfgy & \xrightarrow{g(\epsilon y)} & gy \end{array}$$

Thus, canceling all but the right-hand homotopy, we have $g(\epsilon y) = \eta(gy)$ as desired. \square

However, it is important that we do not include *both* τ and v in the definition of $\text{ishae}(f)$ (whence the name “half adjoint equivalence”). If we did, then after canceling contractible types we would still have one remaining datum — unless we added another higher coherence condition. In general, we expect to get a well-behaved type if we cut off after an odd number of coherences.

Of course, it is obvious that $\text{ishae}(f) \rightarrow \text{qinv}(f)$: simply forget the coherence datum. The other direction is a version of a standard argument from homotopy theory and category theory.

Theorem 4.2.3. *For any $f : A \rightarrow B$ we have $\text{qinv}(f) \rightarrow \text{ishae}(f)$.*

Proof. Suppose that (g, η, ϵ) is a quasi-inverse for f . We have to provide a quadruple $(g', \eta', \epsilon', \tau)$ witnessing that f is a half adjoint equivalence. To define g' and η' , we can just make the obvious choice by setting $g' := g$ and $\eta' := \eta$. However, in the definition of ϵ' we need start worrying about the construction of τ , so we cannot just follow our nose and take ϵ' to be ϵ . Instead, we take

$$\epsilon'(b) := \epsilon(f(g(b)))^{-1} \bullet (f(\eta(g(b))) \bullet \epsilon(b)).$$

Now we need to find

$$\tau(a) : f(\eta(a)) = \epsilon(f(g(f(a))))^{-1} \bullet (f(\eta(g(f(a)))) \bullet \epsilon(f(a))).$$

Note first that by Corollary 2.4.4, we have $\eta(g(f(a))) = g(f(\eta(a)))$. Therefore, we can apply Lemma 2.4.3 to compute

$$\begin{aligned} f(\eta(g(f(a)))) \bullet \epsilon(f(a)) &= f(g(f(\eta(a)))) \bullet \epsilon(f(a)) \\ &= \epsilon(f(g(f(a)))) \bullet f(\eta(a)) \end{aligned}$$

from which we get the desired path $\tau(a)$. \square

Combining this with Lemma 4.2.2 (or symmetrizing the proof), we also have $\text{qinv}(f) \rightarrow \text{ishae}'(f)$.

It remains to show that $\text{ishae}(f)$ is a mere proposition. For this, we will need to know that the fibers of an equivalence are contractible.

Definition 4.2.4. The **fiber** of a map $f : A \rightarrow B$ over a point $y : B$ is

$$\text{fib}_f(y) := \sum_{x:A} (f(x) = y).$$

In homotopy theory, this is what would be called the *homotopy fiber* of f . The path lemmas in §2.5 yield the following characterization of paths in fibers:

Lemma 4.2.5. For any $f : A \rightarrow B$, $y : B$, and $(x, p), (x', p') : \text{fib}_f(y)$, we have

$$((x, p) = (x', p')) \simeq \left(\sum_{\gamma:x=x'} f(\gamma) \bullet p' = p \right)$$

Theorem 4.2.6. If $f : A \rightarrow B$ is a half adjoint equivalence, then for any $y : B$ the fiber $\text{fib}_f(y)$ is contractible.

Proof. Let $(g, \eta, \epsilon, \tau) : \text{ishae}(f)$, and fix $y : B$. As our center of contraction for $\text{fib}_f(y)$ we choose $(gy, \epsilon y)$. Now take any $(x, p) : \text{fib}_f(y)$; we want to construct a path from $(gy, \epsilon y)$ to (x, p) . By Lemma 4.2.5, it suffices to give a path $\gamma : gy = x$ such that $f(\gamma) \bullet p = \epsilon y$. We put $\gamma := g(p)^{-1} \bullet \eta x$. Then we have

$$\begin{aligned} f(\gamma) \bullet p &= fg(p)^{-1} \bullet f(\eta x) \bullet p \\ &= fg(p)^{-1} \bullet \epsilon(fx) \bullet p \\ &= \epsilon y \end{aligned}$$

where the second equality follows by τx and the third equality is naturality of ϵ . \square

We now define the types which encapsulate contractible pairs of data. The following types put together the quasi-inverse g with one of the homotopies.

Definition 4.2.7. Given a function $f : A \rightarrow B$, we define the types

$$\begin{aligned} \text{linv}(f) &:= \sum_{g:B \rightarrow A} (g \circ f \sim \text{id}_A) \\ \text{rinv}(f) &:= \sum_{g:B \rightarrow A} (f \circ g \sim \text{id}_B) \end{aligned}$$

of **left inverses** and **right inverses** to f , respectively. We call f **left invertible** if $\text{linv}(f)$ is inhabited, and similarly **right invertible** if $\text{rinv}(f)$ is inhabited.

Lemma 4.2.8. If $f : A \rightarrow B$ has a quasi-inverse, then so do

$$\begin{aligned} (f \circ -) &: (C \rightarrow A) \rightarrow (C \rightarrow B) \\ (- \circ f) &: (B \rightarrow C) \rightarrow (A \rightarrow C). \end{aligned}$$

Proof. If g is a quasi-inverse of f , then $(g \circ -)$ and $(- \circ g)$ are quasi-inverses of $(f \circ -)$ and $(- \circ f)$ respectively. \square

Lemma 4.2.9. If $f : A \rightarrow B$ has a quasi-inverse, then the types $\text{rinv}(f)$ and $\text{linv}(f)$ are contractible.

Proof. By function extensionality, we have

$$\text{linv}(f) \simeq \sum_{g:B \rightarrow A} (g \circ f = \text{id}_A).$$

But this is the fiber of $(-\circ f)$ over id_A , and so by Lemma 4.2.8 and Theorems 4.2.3 and 4.2.6, it is contractible. Similarly, $\text{rinv}(f)$ is equivalent to the fiber of $(f \circ -)$ over id_B and hence contractible. \square

Next we define the types which put together the other homotopy with the additional coherence datum.

Definition 4.2.10. For $f : A \rightarrow B$, a left inverse $(g, \eta) : \text{linv}(f)$, and a right inverse $(g, \epsilon) : \text{rinv}(f)$, we denote

$$\begin{aligned} \text{lcoh}_f(g, \eta) &:= \sum_{(\epsilon: f \circ g \sim \text{id}_B)} \prod_{(y:B)} g(\epsilon y) = \eta(gy), \\ \text{rcoh}_f(g, \epsilon) &:= \sum_{(\eta: g \circ f \sim \text{id}_A)} \prod_{(x:A)} f(\eta x) = \epsilon(fx). \end{aligned}$$

Lemma 4.2.11. For any f, g, ϵ, η , we have

$$\begin{aligned} \text{lcoh}_f(g, \eta) &\simeq \prod_{y:B} (fgy, \eta(gy)) =_{\text{fib}_g(gy)} (y, \text{refl}_{gy}), \\ \text{rcoh}_f(g, \epsilon) &\simeq \prod_{x:A} (gfx, \epsilon(fx)) =_{\text{fib}_f(fx)} (x, \text{refl}_{fx}). \end{aligned}$$

Proof. Using Lemma 4.2.5. \square

Lemma 4.2.12. If f is a half adjoint equivalence, then for any $(g, \epsilon) : \text{rinv}(f)$, the type $\text{rcoh}_f(g, \epsilon)$ is contractible.

Proof. By Lemma 4.2.11 and the fact that dependent function types preserve contractible spaces, it suffices to show that for each $x : A$, the type $(gfix, \epsilon(fx)) =_{\text{fib}_f(fx)} (x, \text{refl}_{fx})$ is contractible. But by Theorem 4.2.6, $\text{fib}_f(fx)$ is contractible, and any path space of a contractible space is itself contractible. \square

Theorem 4.2.13. For any $f : A \rightarrow B$, the type $\text{ishae}(f)$ is a mere proposition.

Proof. By Exercise 3.5 it suffices to assume f to be a half adjoint equivalence and show that $\text{ishae}(f)$ is contractible. Now by associativity of Σ (Exercise 2.10), the type $\text{ishae}(f)$ is equivalent to

$$\sum_{u:\text{rinv}(f)} \text{rcoh}_f(\text{pr}_1(u), \text{pr}_2(u)).$$

But by Lemmas 4.2.9 and 4.2.12 and the fact that Σ preserves contractibility, the latter type is also contractible. \square

Thus, we have shown that $\text{ishae}(f)$ has all three desiderata for the type $\text{isequiv}(f)$. In the next two sections we consider a couple of other possibilities.

4.3 Bi-invertible maps

Using the language introduced in §4.2, we can restate the definition proposed in §2.4 as follows.

Definition 4.3.1. We say $f : A \rightarrow B$ is **bi-invertible** if it has both a left inverse and a right inverse:

$$\text{biinv}(f) := \text{linv}(f) \times \text{rinv}(f).$$

In §2.4 we proved that $\text{qinv}(f) \rightarrow \text{biinv}(f)$ and $\text{biinv}(f) \rightarrow \text{qinv}(f)$. What remains is the following.

Theorem 4.3.2. *For any $f : A \rightarrow B$, the type $\text{biinv}(f)$ is a mere proposition.*

Proof. We may suppose f to be bi-invertible and show that $\text{biinv}(f)$ is contractible. But since $\text{biinv}(f) \rightarrow \text{qinv}(f)$, by Lemma 4.2.9 in this case both $\text{linv}(f)$ and $\text{rinv}(f)$ are contractible, and the product of contractible types is contractible. \square

Note that this also fits the proposal made at the beginning of §4.2: we combine g and η into a contractible type and add an additional datum which combines with ϵ into a contractible type. The difference is that instead of adding a *higher* datum (a 2-dimensional path) to combine with ϵ , we add a *lower* one (a right inverse that is separate from the left inverse).

Corollary 4.3.3. *For any $f : A \rightarrow B$ we have $\text{biinv}(f) \simeq \text{ishae}(f)$.*

Proof. We have $\text{biinv}(f) \rightarrow \text{qinv}(f) \rightarrow \text{ishae}(f)$ and $\text{ishae}(f) \rightarrow \text{qinv}(f) \rightarrow \text{biinv}(f)$. Since both $\text{ishae}(f)$ and $\text{biinv}(f)$ are mere propositions, the equivalence follows from Lemma 3.3.3. \square

4.4 Contractible fibers

Note that our proofs about $\text{ishae}(f)$ and $\text{biinv}(f)$ made essential use of the fact that the fibers of an equivalence are contractible. In fact, it turns out that this property is itself a sufficient definition of equivalence.

Definition 4.4.1 (Contractible maps). A map $f : A \rightarrow B$ is **contractible** if for all $y : B$, the fiber $\text{fib}_f(y)$ is contractible.

Thus, the type $\text{isContr}(f)$ is defined to be

$$\text{isContr}(f) := \prod_{y:B} \text{isContr}(\text{fib}_f(y)) \tag{4.4.2}$$

Note that in §3.11 we defined what it means for a *type* to be contractible. Here we are defining what it means for a *map* to be contractible. Our terminology follows the general homotopy-theoretic practice of saying that a map has a certain property if all of its (homotopy) fibers have that property. Thus, a type A is contractible just when the map $A \rightarrow \mathbf{1}$ is contractible. From Chapter 7 onwards we will also call contractible maps and types (-2) -truncated.

We have already shown in Theorem 4.2.6 that $\text{ishae}(f) \rightarrow \text{isContr}(f)$. Conversely:

Theorem 4.4.3. *For any $f : A \rightarrow B$ we have $\text{isContr}(f) \rightarrow \text{ishae}(f)$.*

Proof. Let $P : \text{isContr}(f)$. We define an inverse mapping $g : B \rightarrow A$ by sending each $y : B$ to the center of contraction of the fiber at y :

$$g(y) := \text{pr}_1(\text{pr}_1(Py)).$$

We can thus define the homotopy ϵ by mapping y to the witness that $g(y)$ indeed belongs to the fiber at y :

$$\epsilon(y) := \text{pr}_2(\text{pr}_1(Py)).$$

It remains to define η and τ . This of course amounts to giving an element of $\text{rcoh}_f(g, \epsilon)$. By Lemma 4.2.11, this is the same as giving for each $x : A$ a path from $(g(fx), \epsilon(fx))$ to (x, refl_{fx}) in the fiber of f over fx . But this is easy: for any $x : A$, the type $\text{fib}_f(fx)$ is contractible by assumption, hence such a path must exist. We can construct it explicitly as

$$(\text{pr}_2(P(fx))(g(fx), \epsilon(fx)))^{-1} \cdot (\text{pr}_2(P(fx))(x, \text{refl}_{fx})). \quad \square$$

It is also easy to see:

Lemma 4.4.4. *For any f , the type $\text{isContr}(f)$ is a mere proposition.*

Proof. By Lemma 3.11.4, each type $\text{isContr}(\text{fib}_f(y))$ is a mere proposition. Thus, by Example 3.6.2, so is (4.4.2). \square

Theorem 4.4.5. *For any $f : A \rightarrow B$ we have $\text{isContr}(f) \simeq \text{ishae}(f)$.*

Proof. We have already established a logical equivalence $\text{isContr}(f) \Leftrightarrow \text{ishae}(f)$, and both are mere propositions (Lemma 4.4.4 and Theorem 4.2.13). Thus, Lemma 3.3.3 applies. \square

Usually, we prove that a function is an equivalence by exhibiting a quasi-inverse, but sometimes this definition is more convenient. For instance, it implies that when proving a function to be an equivalence, we are free to assume that its codomain is inhabited.

Corollary 4.4.6. *If $f : A \rightarrow B$ is such that $B \rightarrow \text{isequiv}(f)$, then f is an equivalence.*

Proof. To show f is an equivalence, it suffices to show that $\text{fib}_f(y)$ is contractible for any $y : B$. But if $e : B \rightarrow \text{isequiv}(f)$, then given any such y we have $e(y) : \text{isequiv}(f)$, so that f is an equivalence and hence $\text{fib}_f(y)$ is contractible, as desired. \square

4.5 On the definition of equivalences

We have shown that all three definitions of equivalence satisfy the three desirable properties and are pairwise equivalent:

$$\text{isContr}(f) \simeq \text{ishae}(f) \simeq \text{biinv}(f).$$

(There are yet more possible definitions of equivalence, but we will stop with these three. See Exercise 3.8 and the exercises in this chapter for some more.) Thus, we may choose any one of them as “the” definition of $\text{isequiv}(f)$. For definiteness, we choose to define

$$\text{isequiv}(f) := \text{ishae}(f).$$

This choice is advantageous for formalization, since $\text{ishae}(f)$ contains the most directly useful data. On the other hand, for other purposes, $\text{biinv}(f)$ is often easier to deal with, since it contains no 2-dimensional paths and its two symmetrical halves can be treated independently. However, for purposes of this book, the specific choice will make little difference.

In the rest of this chapter, we study some other properties and characterizations of equivalences.

4.6 Surjections and embeddings

When A and B are sets and $f : A \rightarrow B$ is an equivalence, we also call it as **isomorphism** or a **bijection**. (We avoid these words for types that are not sets, since in homotopy theory and higher category theory they often denote a stricter notion of “sameness” than homotopy equivalence.) In set theory, a function is a bijection just when it is both injective and surjective. The same is true in type theory, if we formulate these conditions appropriately. For clarity, when dealing with types that are not sets, we will speak of *embeddings* instead of injections.

Definition 4.6.1. Let $f : A \rightarrow B$.

- (i) We say f is **surjective** (or a **surjection**) if for every $b : B$ we have $\|\text{fib}_f(b)\|$.
- (ii) We say f is an **embedding** if for every $x, y : A$ the function $\text{ap}_f : (x =_A y) \rightarrow (f(x) =_B f(y))$ is an equivalence.

In other words, f is surjective if every fiber of f is merely inhabited, or equivalently if for all $b : B$ there merely exists an $a : A$ such that $f(a) = b$. In traditional logical notation, f is surjective if $\forall(b : B). \exists(a : A). (f(a) = b)$. This must be distinguished from the stronger assertion that $\prod_{(b:B)} \sum_{(a:A)} (f(a) = b)$; if this holds we say that f is a **split surjection**. (Since this latter type is equivalent to $\sum_{(g:B \rightarrow A)} \prod_{(b:B)} (f(g(b)) = b)$, being a split surjection is the same as being a *retraction* as defined in §3.11.)

The axiom of choice from §3.8 says exactly that every surjection *between sets* is split. However, in the presence of the univalence axiom, it is simply false that *all* surjections are split. In Lemma 3.8.5 we constructed a type family $Y : X \rightarrow \mathcal{U}$ such that $\prod_{(x:X)} \|Y(x)\|$ but $\neg \prod_{(x:X)} Y(x)$; for any such family, the first projection $(\sum_{(x:X)} Y(x)) \rightarrow X$ is a surjection that is not split.

If A and B are sets, then by Lemma 3.3.3, f is an embedding just when

$$\prod_{x,y:A} (f(x) =_B f(y)) \rightarrow (x =_A y). \quad (4.6.2)$$

In this case we say that f is **injective**, or an **injection**. We avoid these word for types that are not sets, because they might be interpreted as (4.6.2), which is an ill-behaved notion for non-sets. It is also true that any function between sets is surjective if and only if it is an *epimorphism* in a suitable sense, but this also fails for more general types, and surjectivity is generally the more important notion.

Theorem 4.6.3. *A function $f : A \rightarrow B$ is an equivalence if and only if it is both surjective and an embedding.*

Proof. If f is an equivalence, then each $\text{fib}_f(b)$ is contractible, hence so is $\|\text{fib}_f(b)\|$, so f is surjective. And we showed in Theorem 2.11.1 that any equivalence is an embedding.

Conversely, suppose f is a surjective embedding. Let $b : B$; we show that $\sum_{(x:A)} (f(x) = b)$ is contractible. Since f is surjective, there merely exists an $a : A$ such that $f(a) = b$. Thus, the fiber of f over b is inhabited; it remains to show it is a mere proposition. For this, suppose given $x, y : A$ with $p : f(x) = b$ and $q : f(y) = b$. Then since ap_f is an equivalence, there exists $r : x = y$ with $\text{ap}_f(r) = p \bullet q^{-1}$. However, using the characterization of paths in Σ -types, the latter equality rearranges to $r_*(p) = q$. Thus, together with r it exhibits $(x, p) = (y, q)$ in the fiber of f over b . \square

Corollary 4.6.4. *For any $f : A \rightarrow B$ we have*

$$\text{isequiv}(f) \simeq (\text{isEmbedding}(f) \times \text{isSurjective}(f)).$$

Proof. Being a surjection and an embedding are both mere propositions; now apply Lemma 3.3.3. \square

Of course, this cannot be used as a definition of “equivalence”, since the definition of embeddings refers to equivalences. However, this characterization can still be useful; see ???. We will generalize it in Chapter 7.

4.7 Closure properties of equivalences

We have already seen in Lemma 2.4.12 that equivalences are closed under composition. Furthermore, we have:

Theorem 4.7.1 (The 2-out-of-3 property). *Suppose $f : A \rightarrow B$ and $g : B \rightarrow C$. If any two of f , g , and $g \circ f$ are equivalences, so is the third.*

Proof. If $g \circ f$ and g are equivalences, then $(g \circ f)^{-1} \circ g$ is a quasi-inverse to f . On the one hand, we have $(g \circ f)^{-1} \circ g \circ f \sim \text{id}_A$, while on the other we have

$$\begin{aligned} f \circ (g \circ f)^{-1} \circ g &\sim g^{-1} \circ g \circ f \circ (g \circ f)^{-1} \circ g \\ &\sim g^{-1} \circ g \\ &\sim \text{id}_B. \end{aligned}$$

Similarly, if $g \circ f$ and f are equivalences, then $f \circ (g \circ f)^{-1}$ is a quasi-inverse to g . \square

This is a standard closure condition on equivalences from homotopy theory. Also well-known is that they are closed under retracts, in the following sense.

Definition 4.7.2. A function $g : A \rightarrow B$ is said to be a **retract** of a function $f : X \rightarrow Y$ if there is a diagram

$$\begin{array}{ccccc} A & \xrightarrow{s} & X & \xrightarrow{r} & A \\ g \downarrow & & f \downarrow & & \downarrow g \\ B & \xrightarrow{s'} & Y & \xrightarrow{r'} & B \end{array}$$

for which there are

- (i) a homotopy $R : r \circ s \sim \text{id}_A$.
- (ii) a homotopy $R' : r' \circ s' \sim \text{id}_B$.
- (iii) a homotopy $L : f \circ s \sim s' \circ g$.
- (iv) a homotopy $K : g \circ r \sim r' \circ f$.
- (v) for every $a : A$, a path $H(a)$ witnessing the commutativity of the square

$$\begin{array}{ccc} g(r(s(a))) & \xrightarrow{K(s(a))} & r'(f(s(a))) \\ g(R(a)) \parallel & & \parallel r'(L(a)) \\ g(a) & \xrightarrow{R'(g(a))^{-1}} & r'(s'(g(a))) \end{array}$$

Recall that in §3.11 we defined what it means for a type to be a retract of another. This is a special case of the above definition where B and Y are $\mathbf{1}$. Conversely, just as with contractibility, retractions of maps induce retractions of their fibers.

Lemma 4.7.3. *If a function $g : A \rightarrow B$ is a retract of a function $f : X \rightarrow Y$, then $\text{fib}_g(b)$ is a retract of $\text{fib}_f(s'(b))$ for every $b : B$, where $s' : B \rightarrow Y$ is as in Definition 4.7.2.*

Proof. Suppose that $g : A \rightarrow B$ is a retract of $f : X \rightarrow Y$. Then for any $b : B$ we have the functions

$$\begin{aligned}\varphi_b : \text{fib}_g(b) &\rightarrow \text{fib}_f(s'(b)), & \varphi_b(a, p) &:= (s(a), L(a) \bullet s'(p)), \\ \psi_b : \text{fib}_f(s'(b)) &\rightarrow \text{fib}_g(b), & \psi_b(x, q) &:= (r(x), K(x) \bullet r'(q) \bullet R'(b)).\end{aligned}$$

Then we have $\psi_b(\varphi_b(a, p)) \equiv (r(s(a)), K(s(a)) \bullet r'(L(a) \bullet s'(p)) \bullet R'(b))$. We claim ψ_b is a retraction with section φ_b for all $b : B$, which is to say that for all $(a, p) : \text{fib}_g(b)$ we have $\psi_b(\varphi_b(a, p)) = (a, p)$. In other words, we want to show

$$\prod_{(b:B)} \prod_{(a:A)} \prod_{(p:g(a)=b)} \psi_b(\varphi_b(a, p)) = (a, p).$$

By reordering the first two Π s and applying a version of Lemma 3.11.9, this is equivalent to

$$\prod_{a:A} \psi_{g(a)}(\varphi_{g(a)}(a, \text{refl}_{g(a)})) = (a, \text{refl}_{g(a)}).$$

For any a , by Theorem 2.7.2, this equality of pairs is equivalent to a pair of equalities. The first components are equal by $R(a) : r(s(a)) = a$, so we need only show

$$R(a)_*(K(s(a)) \bullet r'(L(a)) \bullet R'(g(a))) = \text{refl}_{g(a)}.$$

But this transportation computes as $g(R(a))^{-1} \bullet K(s(a)) \bullet r'(L(a)) \bullet R'(g(a))$, so the required path is given by $H(a)$. \square

Theorem 4.7.4. *If g is a retract of an equivalence f , then g is also an equivalence.*

Proof. By Lemma 4.7.3, every fiber of g is a retract of a fiber of f . Thus, by Lemma 3.11.7, if the latter are all contractible, so are the former. \square

Finally, we show that fiberwise equivalences can be characterized in terms of equivalences of total spaces. To explain the terminology, recall from §2.3 that a type family $P : A \rightarrow \mathcal{U}$ can be viewed as a fibration over A with total space $\sum_{(x:A)} P(x)$, the fibration being the projection $\text{pr}_1 : \sum_{(x:A)} P(x) \rightarrow A$. From this point of view, given two type families $P, Q : A \rightarrow \mathcal{U}$, we may refer to a function $f : \prod_{(x:A)} (P(x) \rightarrow Q(x))$ as a **fiberwise map** or a **fiberwise transformation**. Such a map induces a function on total spaces:

Definition 4.7.5. Given type families $P, Q : A \rightarrow \mathcal{U}$ and a map $f : \prod_{(x:A)} P(x) \rightarrow Q(x)$, we define

$$\text{total}(f) := \lambda w. (\text{pr}_1 w, f(\text{pr}_1 w, \text{pr}_2 w)) : \sum_{x:A} P(x) \rightarrow \sum_{x:A} Q(x).$$

Theorem 4.7.6. *Suppose that f is a fiberwise transformation between families P and Q over a type A and let $x : A$ and $v : Q(x)$. Then we have an equivalence*

$$\text{fib}_{\text{total}(f)}((x, v)) \simeq \text{fib}_{f(x)}(v).$$

Proof. We calculate:

$$\begin{aligned}
\mathbf{fib}_{\text{total}(f)}((x, v)) &\equiv \sum_{w:\sum_{(x:A)} P(x)} (\mathbf{pr}_1 w, f(\mathbf{pr}_1 w, \mathbf{pr}_2 w)) = (x, v) \\
&\simeq \sum_{(a:A)} \sum_{(u:P(a))} (a, f(a, u)) = (x, v) && \text{(by Exercise 2.10)} \\
&\simeq \sum_{(a:A)} \sum_{(u:P(a))} \sum_{(p:a=x)} p_*(f(a, u)) = v && \text{(by Theorem 2.7.2)} \\
&\simeq \sum_{(a:A)} \sum_{(p:a=x)} \sum_{(u:P(a))} p_*(f(a, u)) = v \\
&\simeq \sum_{u:P(x)} f(x, u) = v \\
&\equiv \mathbf{fib}_{f(x)}(v).
\end{aligned} \tag{*}$$

The equivalence (*) follows from Lemmas 3.11.8 and 3.11.9 and Exercise 2.10. \square

We say that a fiberwise transformation $f : \prod_{(x:A)} P(x) \rightarrow Q(x)$ is a **fiberwise equivalence** if each $f(x) : P(x) \rightarrow Q(x)$ is an equivalence.

Theorem 4.7.7. Suppose that f is a fiberwise transformation between families P and Q over a type A . Then f is a fiberwise equivalence if and only if $\text{total}(f)$ is an equivalence.

Proof. Let f, P, Q and A be as in the statement of the theorem. By Theorem 4.7.6 it follows for all $x : A$ and $v : Q(x)$ that $\mathbf{fib}_{\text{total}(f)}((x, v))$ is contractible if and only if $\mathbf{fib}_{f(x)}(v)$ is contractible. Thus, $\mathbf{fib}_{\text{total}(f)}(w)$ is contractible for all $w : \sum_{(x:A)} Q(x)$ if and only if $\mathbf{fib}_{f(x)}(v)$ is contractible for all $x : A$ and $v : Q(x)$. \square

4.8 The object classifier

In type theory we have a basic notion of *family of types*, namely a function $B : A \rightarrow \mathcal{U}$. We have seen that such families behave somewhat like *fibrations* in homotopy theory, with the fibration being the projection $\mathbf{pr}_1 : \sum_{(a:A)} B(a) \rightarrow A$. A basic fact in homotopy theory is that every map is equivalent to a fibration. With univalence at our disposal, we can prove the same thing in type theory.

Lemma 4.8.1. For any type family $B : A \rightarrow \mathcal{U}$, the fiber of $\mathbf{pr}_1 : \sum_{(x:A)} B(x) \rightarrow A$ over $a : A$ is equivalent to $B(a)$:

$$\mathbf{fib}_{\mathbf{pr}_1}(a) \simeq B(a)$$

Proof. We have

$$\begin{aligned}
\mathbf{fib}_{\mathbf{pr}_1}(a) &:= \sum_{u:\sum_{(x:A)} B(x)} \mathbf{pr}_1(u) = a \\
&\simeq \sum_{(x:A)} \sum_{(b:B(x))} (x = a) \\
&\simeq \sum_{(x:A)} \sum_{(p:x=a)} B(x) \\
&\simeq B(a)
\end{aligned}$$

using the left universal property of identity types. \square

Lemma 4.8.2. For any function $f : A \rightarrow B$, we have $A \simeq \sum_{(b:B)} \text{fib}_f(b)$.

Proof. We have

$$\begin{aligned} \sum_{b:B} \text{fib}_f(b) &:= \sum_{(b:B)} \sum_{(a:A)} (f(a) = b) \\ &\simeq \sum_{(a:A)} \sum_{(b:B)} (f(a) = b) \\ &\simeq A \end{aligned}$$

using the fact that $\sum_{(b:B)} (f(a) = b)$ is contractible. \square

Theorem 4.8.3. For any type B there is an equivalence

$$\chi : \left(\sum_{A:\mathcal{U}} (A \rightarrow B) \right) \simeq (B \rightarrow \mathcal{U}).$$

Proof. We have to construct quasi-inverses

$$\begin{aligned} \chi : \left(\sum_{A:\mathcal{U}} (A \rightarrow B) \right) &\rightarrow B \rightarrow \mathcal{U} \\ \psi : (B \rightarrow \mathcal{U}) &\rightarrow \left(\sum_{A:\mathcal{U}} (A \rightarrow B) \right). \end{aligned}$$

We define χ by $\chi((A, f), b) := \text{fib}_f(b)$, and ψ by $\psi(P) := ((\sum_{(b:B)} P(b)), \text{pr}_1)$. Now we have to verify that $\chi \circ \psi \sim \text{id}$ and that $\psi \circ \chi \sim \text{id}$.

- (i) Let $P : B \rightarrow \mathcal{U}$. By Lemma 4.8.1, $\text{fib}_{\text{pr}_1}(b) \simeq P(b)$ for any $b : B$, so it follows immediately that $P \sim \chi(\psi(P))$.
- (ii) Let $f : A \rightarrow B$ be a function. We have to find a path

$$\left(\sum_{(b:B)} \text{fib}_f(b), \text{pr}_1 \right) = (A, f).$$

First note that by Lemma 4.8.2, we have $e : \sum_{(b:B)} \text{fib}_f(b) \simeq A$ with $e(b, a, p) := a$ and $e^{-1}(a) := (f(a), a, \text{refl}_{f(a)})$. By Theorem 2.7.2, it remains to show $(\text{ua}(e))_*(\text{pr}_1) = f$. But by the computation rule for univalence and (2.9.4), we have $(\text{ua}(e))_*(\text{pr}_1) = \text{pr}_1 \circ e^{-1}$, and the definition of e^{-1} immediately yields $\text{pr}_1 \circ e^{-1} \equiv f$. \square

In particular, this implies that we have an *object classifier* in the sense of higher topos theory. Recall from Definition 2.1.7 that \mathcal{U}_\bullet denotes the type $\sum_{(A:\mathcal{U})} A$ of pointed types.

Theorem 4.8.4. Let $f : A \rightarrow B$ be a function. Then the diagram

$$\begin{array}{ccc} A & \xrightarrow{\vartheta_f} & \mathcal{U}_\bullet \\ f \downarrow & & \downarrow \text{pr}_1 \\ B & \xrightarrow{\chi_f} & \mathcal{U} \end{array}$$

is a pullback square (see Exercise 2.11). Here the function ϑ_f is defined by

$$\lambda a. (\text{fib}_f(f(a)), (a, \text{refl}_{f(a)})).$$

Proof. Note that we have the equivalences

$$\begin{aligned} A &\simeq \sum_{b:B} \text{fib}_f(b) \\ &\simeq \sum_{(b:B)} \sum_{(X:\mathcal{U})} \sum_{(p:\text{fib}_f(b)=X)} X \\ &\simeq \sum_{(b:B)} \sum_{(X:\mathcal{U})} \sum_{(x:X)} \text{fib}_f(b) = X \\ &\simeq \sum_{(b:B)} \sum_{(Y:\mathcal{U}_\bullet)} \text{fib}_f(b) = \text{pr}_1 Y \\ &\equiv B \times_{\mathcal{U}} \mathcal{U}_\bullet \end{aligned}$$

which gives us a composite equivalence $e : A \simeq B \times_{\mathcal{U}} \mathcal{U}_\bullet$. We may display the action of this composite equivalence step by step by

$$\begin{aligned} a &\mapsto (f(a), (a, \text{refl}_{f(a)})) \\ &\mapsto (f(a), \text{fib}_f(f(a)), \text{refl}_{\text{fib}_f(f(a))}, (a, \text{refl}_{f(a)})) \\ &\mapsto (f(a), \text{fib}_f(f(a)), (a, \text{refl}_{f(a)}), \text{refl}_{\text{fib}_f(f(a))}). \end{aligned}$$

Therefore, we get homotopies $f \sim \text{pr}_1 \circ e$ and $\vartheta_f \sim \text{pr}_2 \circ e$. \square

4.9 Univalence implies function extensionality

In the last section of this chapter we include a proof that the univalence axiom implies function extensionality. Thus, in this section we work *without* the function extensionality axiom. The proof consists of two steps. First we show in Theorem 4.9.4 that the univalence axiom implies a weak form of function extensionality, defined in Definition 4.9.1 below. The principle of weak function extensionality in turn implies the usual function extensionality, and it does so without the univalence axiom (Theorem 4.9.5).

Let \mathcal{U} be a universe; we will explicitly indicate where we assume that it is univalent.

Definition 4.9.1. The **weak function extensionality principle** asserts that there is a function

$$\left(\prod_{x:A} \text{isContr}(P(x)) \right) \rightarrow \text{isContr}\left(\prod_{x:A} P(x) \right)$$

for any family $P : A \rightarrow \mathcal{U}$ of types over any type A .

The following lemma is easy to prove using function extensionality; the point here is that it also follows from univalence without assuming function extensionality separately.

Lemma 4.9.2. *Assuming \mathcal{U} is univalent, for any $A, B, X : \mathcal{U}$ and any $e : A \simeq B$, there is an equivalence*

$$(X \rightarrow A) \simeq (X \rightarrow B)$$

of which the underlying map is given by post-composition with the underlying function of e .

Proof. As in the proof of Lemma 4.1.1, we may assume that $e = \text{idtoeqv}(p)$ for some $p : A = B$. Then by path induction, we may assume p is refl_A , so that $e = \text{id}_A$. But in this case, post-composition with e is the identity, hence an equivalence. \square

Corollary 4.9.3. *Let $P : A \rightarrow \mathcal{U}$ be a family of contractible types, i.e. $\prod_{(x:A)} \text{isContr}(P(x))$. Then the projection $\text{pr}_1 : (\sum_{(x:A)} P(x)) \rightarrow A$ is an equivalence. Assuming \mathcal{U} is univalent, it follows immediately that post-composition with pr_1 gives an equivalence*

$$\alpha : \left(A \rightarrow \sum_{x:A} P(x) \right) \simeq (A \rightarrow A).$$

Proof. By Lemma 4.8.1, for $\text{pr}_1 : \sum_{(x:A)} P(x) \rightarrow A$ and $x : A$ we have an equivalence

$$\text{fib}_{\text{pr}_1}(x) \simeq P(x).$$

Therefore pr_1 is an equivalence whenever each $P(x)$ is contractible. The assertion is now a consequence of Lemma 4.9.2. \square

In particular, the homotopy fiber of the above equivalence at id_A is contractible. Therefore, we can show that univalence implies weak function extensionality by showing that the dependent function type $\prod_{(x:A)} P(x)$ is a retract of $\text{fib}_\alpha(\text{id}_A)$.

Theorem 4.9.4. *In a univalent universe \mathcal{U} , suppose that $P : A \rightarrow \mathcal{U}$ is a family of contractible types and let α be the function of Corollary 4.9.3. Then $\prod_{(x:A)} P(x)$ is a retract of $\text{fib}_\alpha(\text{id}_A)$. As a consequence, $\prod_{(x:A)} P(x)$ is contractible. In other words, the univalence axiom implies the weak function extensionality principle.*

Proof. Define the functions

$$\begin{aligned} \varphi : & (\prod_{(x:A)} P(x)) \rightarrow \text{fib}_\alpha(\text{id}_A), \\ \varphi(f) : & \equiv (\lambda x. (x, f(x)), \text{refl}_{\text{id}_A}), \end{aligned}$$

and

$$\begin{aligned} \psi : & \text{fib}_\alpha(\text{id}_A) \rightarrow \prod_{(x:A)} P(x), \\ \psi(g, p) : & \equiv \lambda x. \text{happly}(p, x)_*(\text{pr}_2(g(x))). \end{aligned}$$

Then $\psi(\varphi(f)) = \lambda x. f(x)$, which is f , by the uniqueness principle for dependent function types. \square

We now show that weak function extensionality implies the usual function extensionality. Recall from (2.9.2) the function $\text{happly}(f, g) : (f = g) \rightarrow (f \sim g)$ which converts equality of functions to homotopy. In the proof that follows, the univalence axiom is not used.

Theorem 4.9.5. *Weak function extensionality implies the function extensionality Axiom 2.9.3.*

Proof. We want to show that

$$\prod_{(A:\mathcal{U})} \prod_{(P:A \rightarrow \mathcal{U})} \prod_{(f,g:\prod_{(x:A)} P(x))} \text{isequiv}(\text{happly}(f, g)).$$

Since a fiberwise map induces an equivalence on total spaces if and only if it is fiberwise an equivalence by Theorem 4.7.7, it suffices to show that the function of type

$$\left(\sum_{g:\prod_{(x:A)} P(x)} (f = g) \right) \rightarrow \sum_{g:\prod_{(x:A)} P(x)} (f \sim g)$$

induced by $\lambda(g : \prod_{(x:A)} P(x)) . \text{happly}(f, g)$ is an equivalence. Since the type on the left is contractible by Lemma 3.11.8, it suffices to show that the type on the right:

$$\sum_{(g:\prod_{(x:A)} P(x))} \prod_{(x:A)} f(x) = g(x) \quad (4.9.6)$$

is contractible. Now Theorem 2.15.7 says that this is equivalent to

$$\prod_{(x:A)} \sum_{(u:P(x))} f(x) = u. \quad (4.9.7)$$

The proof of Theorem 2.15.7 uses function extensionality, but only for one of the composites. Thus, without assuming function extensionality, we can conclude that (4.9.6) is a retract of (4.9.7). And (4.9.7) is a product of contractible types, which is contractible by the weak function extensionality principle; hence (4.9.6) is also contractible. \square

Notes

The fact that the space of continuous maps equipped with quasi-inverses has the wrong homotopy type to be the “space of homotopy equivalences” is well-known in algebraic topology. In that context, the “space of homotopy equivalences” ($A \simeq B$) is usually defined simply as the subspace of the function space ($A \rightarrow B$) consisting of the functions that are homotopy equivalences. In type theory, this would correspond most closely to $\sum_{(f:A \rightarrow B)} \|\text{qinv}(f)\|$; see Exercise 3.8.

The first definition of equivalence given in homotopy type theory was the one that we have called $\text{isContr}(f)$, which was due to Voevodsky. The possibility of the other definitions was subsequently observed by various people. The basic theorems about adjoint equivalences such as Lemma 4.2.2 and Theorem 4.2.3 are adaptations of standard facts in higher category theory and homotopy theory. Using bi-invertibility as a definition of equivalences was suggested by André Joyal.

The properties of equivalences discussed in §§4.6 and 4.7 are well-known in homotopy theory. Most of them were first proven in type theory by Voevodsky.

The fact that every function is equivalent to a fibration is a standard fact in homotopy theory. The notion of object classifier in $(\infty, 1)$ -category theory (the categorical analogue of Theorem 4.8.3) is due to Rezk (see [Rez05, Lur09]).

Finally, the fact that univalence implies function extensionality (§4.9) is due to Voevodsky. Our proof is a simplification of his. Exercise 4.9 is also due to Voevodsky.

Exercises

Exercise 4.1. Consider the type of “two-sided adjoint equivalence data” for $f : A \rightarrow B$,

$$\sum_{(g:B \rightarrow A)} \sum_{(\eta:g \circ f \sim \text{id}_A)} \sum_{(\epsilon:f \circ g \sim \text{id}_B)} \left(\prod_{x:A} f(\eta x) = \epsilon(fx) \right) \times \left(\prod_{y:B} g(\epsilon y) = \eta(gy) \right).$$

By Lemma 4.2.2, we know that if f is an equivalence, then this type is inhabited. Give a characterization of this type analogous to Lemma 4.1.1.

Can you give an example showing that this type is not generally a mere proposition? (This will be easier after Chapter 6.)

Exercise 4.2. Show that for any $A, B : \mathcal{U}$, the following type is equivalent to $A \simeq B$.

$$\sum_{R:A \rightarrow B \rightarrow \mathcal{U}} \left(\prod_{a:A} \text{isContr} \left(\sum_{b:B} R(a, b) \right) \right) \times \left(\prod_{b:B} \text{isContr} \left(\sum_{a:A} R(a, b) \right) \right).$$

Can you extract from this a definition of a type satisfying the three desiderata of $\text{isequiv}(f)$?

Exercise 4.3. Reformulate the proof of Lemma 4.1.1 without using univalence.

Exercise 4.4 (The unstable octahedral axiom). Suppose $f : A \rightarrow B$ and $g : B \rightarrow C$ and $b : B$.

- (i) Show that there is a natural map $\text{fib}_{g \circ f}(g(b)) \rightarrow \text{fib}_g(g(b))$ whose fiber over $(b, \text{refl}_{g(b)})$ is equivalent to $\text{fib}_f(b)$.
- (ii) Show that $\text{fib}_{g \circ f}(c) \simeq \sum_{(w:\text{fib}_g(c))} \text{fib}_f(\text{pr}_1 w)$.

Exercise 4.5. Prove that equivalences satisfy the *2-out-of-6 property*: given $f : A \rightarrow B$ and $g : B \rightarrow C$ and $h : C \rightarrow D$, if $g \circ f$ and $h \circ g$ are equivalences, so are f, g, h , and $h \circ g \circ f$. Use this to give a higher-level proof of Theorem 2.11.1.

Exercise 4.6. For $A, B : \mathcal{U}$, define

$$\text{idtoqinv}_{A,B} : (A = B) \rightarrow \sum_{f:A \rightarrow B} \text{qinv}(f)$$

by path induction in the obvious way. Let **qinv-univalence** denote the modified form of the univalence axiom which asserts that for all $A, B : \mathcal{U}$ the function $\text{idtoqinv}_{A,B}$ has a quasi-inverse.

- (i) Show that qinv-univalence can be used instead of univalence in the proof of function extensionality in §4.9.
- (ii) Show that qinv-univalence can be used instead of univalence in the proof of Theorem 4.1.3.
- (iii) Show that qinv-univalence is inconsistent (i.e. allows construction of an inhabitant of $\mathbf{0}$). Thus, the use of a “good” version of isequiv is essential in the statement of univalence.

Exercise 4.7. Show that a function $f : A \rightarrow B$ is an embedding if and only if the following two conditions hold:

- (i) f is *left cancellable*, i.e. for any $x, y : A$, if $f(x) = f(y)$ then $x = y$.
- (ii) For any $x : A$, the map $\text{ap}_f : \Omega(A, x) \rightarrow \Omega(B, f(x))$ is an equivalence.

(In particular, if A is a set, then f is an embedding if and only if it is left-cancellable and $\Omega(B, f(x))$ is contractible for all $x : A$.) Give examples to show that neither of (i) or (ii) implies the other.

Exercise 4.8. Show that the type of left-cancellable functions $\mathbf{2} \rightarrow B$ (see Exercise 4.7) is equivalent to $\sum_{(x,y:B)} (x \neq y)$. Give a similar explicit characterization of the type of embeddings $\mathbf{2} \rightarrow B$.

Exercise 4.9. The **naïve non-dependent function extensionality axiom** says that for $A, B : \mathcal{U}$ and $f, g : A \rightarrow B$ there is a function $(\prod_{(x:A)} f(x) = g(x)) \rightarrow (f = g)$. Modify the argument of §4.9 to show that this axiom implies the full function extensionality axiom (Axiom 2.9.3).

Chapter 5

Induction

In Chapter 1, we introduced many ways to form new types from old ones. Except for (dependent) function types and universes, all these rules are special cases of the general notion of *inductive definition*. In this chapter we study inductive definitions more generally.

5.1 Introduction to inductive types

An *inductive type* X can be intuitively understood as a type “freely generated” by a certain finite collection of *constructors*, each of which is a function (of some number of arguments) with codomain X . This includes functions of zero arguments, which are simply elements of X .

When describing a particular inductive type, we list the constructors with bullets. For instance, the type **2** from §1.8 is inductively generated by the following constructors:

- $0_2 : \mathbf{2}$
- $1_2 : \mathbf{2}$

Similarly, **1** is inductively generated by the constructor:

- $\star : \mathbf{1}$

while **0** is inductively generated by no constructors at all. An example where the constructor functions take arguments is the coproduct $A + B$, which is generated by the two constructors

- $\text{inl} : A \rightarrow A + B$
- $\text{inr} : B \rightarrow A + B$.

And an example with a constructor taking multiple arguments is the cartesian product $A \times B$, which is generated by one constructor

- $(-, -) : A \rightarrow B \rightarrow A \times B$.

Crucially, we also allow constructors of inductive types that take arguments from the inductive type being defined. For instance, the type \mathbb{N} of natural numbers has constructors

- $0 : \mathbb{N}$
- $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$.

Another useful example is the type $\text{List}(A)$ of finite lists of elements of some type A , which has constructors

- $\text{nil} : \text{List}(A)$
- $\text{cons} : A \rightarrow \text{List}(A) \rightarrow \text{List}(A)$.

Intuitively, we should understand an inductive type as being *freely generated* by its constructors. That is, the elements of an inductive type are exactly what can be obtained by starting from nothing and applying the constructors repeatedly. (We will see in §5.8 and Chapter 6 that this conception has to be modified slightly for more general kinds of inductive definitions, but for now it is sufficient.) For instance, in the case of $\mathbf{2}$, we should expect that the only elements are 0_2 and 1_2 . Similarly, in the case of \mathbb{N} , we should expect that every element is either 0 or obtained by applying succ to some “previously constructed” natural number.

Rather than assert properties such as this directly, however, we express them by means of an *induction principle*, also called a (*dependent*) *elimination rule*. We have seen these principles already in Chapter 1. For instance, the induction principle for $\mathbf{2}$ is:

- When proving a statement $E : \mathbf{2} \rightarrow \mathcal{U}$ about *all* inhabitants of $\mathbf{2}$, it suffices to prove it for 0_2 and 1_2 , i.e., to give proofs $e_0 : E(0_2)$ and $e_1 : E(1_2)$.

Furthermore, the resulting proof $\text{ind}_2(E, e_0, e_1) : \prod_{(b:\mathbf{2})} E(b)$ behaves as expected when applied to the constructors 0_2 and 1_2 ; this principle is expressed by the *computation rules*:

- We have $\text{ind}_2(E, e_0, e_1, 0_2) \equiv e_0$.
- We have $\text{ind}_2(E, e_0, e_1, 1_2) \equiv e_1$.

Thus, the induction principle for the type $\mathbf{2}$ of booleans allows us to reason by *case analysis*. Since neither of the two constructors takes any arguments, this is all we need for booleans.

For natural numbers, however, case analysis is generally not sufficient: in the case corresponding to the inductive step $\text{succ}(n)$, we also want to presume that the statement being proven has already been shown for n . This gives us the following induction principle:

- When proving a statement $E : \mathbb{N} \rightarrow \mathcal{U}$ about *all* natural numbers, it suffices to prove it for 0 and for $\text{succ}(n)$, assuming it holds for n , i.e., we construct $e_z : E(0)$ and $e_s : \prod_{(n:\mathbb{N})} E(n) \rightarrow E(\text{succ}(n))$.

As in the case of booleans, we also have the associated computation rules for the function $\text{ind}_{\mathbb{N}}(E, e_z, e_s) : \prod_{(x:\mathbb{N})} E(x)$:

- $\text{ind}_{\mathbb{N}}(E, e_z, e_s, 0) \equiv e_z$.
- $\text{ind}_{\mathbb{N}}(E, e_z, e_s, \text{succ}(n)) \equiv e_s(n, \text{ind}_{\mathbb{N}}(E, e_z, e_s, n))$ for any $n : \mathbb{N}$.

The dependent function $\text{ind}_{\mathbb{N}}(E, e_z, e_s)$ can thus be understood as being defined recursively on the argument $x : \mathbb{N}$, via the functions e_z and e_s which we call the **recurrences**. When x is zero, the function simply returns e_z . When x is the successor of another natural number n , the result is obtained by taking the recurrence e_s and substituting the specific predecessor n and the recursive call value $\text{ind}_{\mathbb{N}}(E, e_z, e_s, n)$.

The induction principles for all the examples mentioned above share this family resemblance. In §5.6 we will discuss a general notion of “inductive definition” and how to derive an appropriate *induction principle* for it, but first we investigate various commonalities between inductive definitions.

For instance, we have remarked in every case in Chapter 1 that from the induction principle we can derive a *recursion principle* in which the codomain is a simple type (rather than a family).

Both induction and recursion principles may seem odd, since they yield only the *existence* of a function without seeming to characterize it uniquely. However, in fact the induction principle is strong enough also to prove its own *uniqueness principle*, as in the following theorem.

Theorem 5.1.1. *Let $f, g : \prod_{(x:\mathbb{N})} E(x)$ be two functions which satisfy the recurrences*

$$e_z : E(0) \quad \text{and} \quad e_s : \prod_{n:\mathbb{N}} E(n) \rightarrow E(\text{succ}(n))$$

up to propositional equality, i.e., such that

$$f(0) = e_z \quad \text{and} \quad g(0) = e_z$$

as well as

$$\begin{aligned} \prod_{n:\mathbb{N}} f(\text{succ}(n)) &= e_s(n, f(n)), \\ \prod_{n:\mathbb{N}} g(\text{succ}(n)) &= e_s(n, g(n)). \end{aligned}$$

Then f and g are equal.

Proof. We use induction on the type family $D(x) := f(x) = g(x)$. For the base case, we have

$$f(0) = e_z = g(0).$$

For the inductive case, assume $n : \mathbb{N}$ such that $f(n) = g(n)$. Then

$$f(\text{succ}(n)) = e_s(n, f(n)) = e_s(n, g(n)) = g(\text{succ}(n)).$$

The first and last equality follow from the assumptions on f and g . The middle equality follows from the inductive hypothesis and the fact that application preserves equality. This gives us pointwise equality between f and g ; invoking function extensionality finishes the proof. \square

Note that the uniqueness principle applies even to functions that only satisfy the recurrences *up to propositional equality*, i.e. a path. Of course, the particular function obtained from the induction principle satisfies these recurrences judgmentally; we will return to this point in §5.5. On the other hand, the theorem itself only asserts a propositional equality between functions (see also Exercise 5.2). From a homotopical viewpoint it is natural to ask whether this path is *coherent*, i.e. whether the equality $f = g$ is unique up to higher paths; in §5.4 we will see that this is in fact the case.

Of course, similar uniqueness theorems for functions can generally be formulated and shown for other inductive types as well. In the next section, we show how this uniqueness property, together with univalence, implies that an inductive type such as the natural numbers is completely characterized by its introduction, elimination, and computation rules.

5.2 Uniqueness of inductive types

We have defined “the” natural numbers to be a particular type \mathbb{N} with particular inductive generators 0 and succ . However, by the general principle of inductive definitions in type theory described in the previous section, there is nothing preventing us from defining *another* type in an identical way. That is, suppose we let \mathbb{N}' be the inductive type generated by the constructors

- $0' : \mathbb{N}'$
- $\text{succ}' : \mathbb{N}' \rightarrow \mathbb{N}'$.

Then \mathbb{N}' will have identical-looking induction and recursion principles to \mathbb{N} . When proving a statement $E : \mathbb{N}' \rightarrow \mathcal{U}$ for all of these “new” natural numbers, it suffices to give the proofs $e_z : E(0')$ and $e_s : \prod_{(n:\mathbb{N}')} E(n) \rightarrow E(\text{succ}'(n))$. And the function $\text{rec}_{\mathbb{N}'}(E, e_z, e_s) : \prod_{(n:\mathbb{N}')} E(n)$ has the following computation rules:

- $\text{rec}_{\mathbb{N}'}(E, e_z, e_s, 0') \equiv e_z$,
- $\text{rec}_{\mathbb{N}'}(E, e_z, e_s, \text{succ}'(n)) \equiv e_s(n, \text{rec}_{\mathbb{N}'}(E, e_z, e_s, n))$ for any $n : \mathbb{N}'$.

But what is the relation between \mathbb{N} and \mathbb{N}' ?

This is not just an academic question, since structures that “look like” the natural numbers can be found in many other places. For instance, we may identify natural numbers with lists over the type with one element (this is arguably the oldest appearance, found on walls of caves), with the non-negative integers, with subsets of the rationals and the reals, and so on. And from a programming point of view, the “unary” representation of our natural numbers is very inefficient, so we might prefer sometimes to use a binary one instead. We would like to be able to identify all of these versions of “the natural numbers” with each other, in order to transfer constructions and results from one to another.

Of course, if two versions of the natural numbers satisfy identical induction principles, then they have identical induced structure. For instance, recall the example of the function `double` defined in §1.9. A similar function for our new natural numbers is readily defined by duplication and adding primes:

$$\text{double}' := \text{rec}_{\mathbb{N}'}(\mathbb{N}', 0', \lambda n. \lambda m. \text{succ}'(\text{succ}'(m))).$$

Simple as this may seem, it has the obvious drawback of leading to a proliferation of duplicates. Not only functions have to be duplicated, but also all lemmas and their proofs. For example, an easy result such as $\prod_{(n:\mathbb{N})} \text{double}(\text{succ}(n)) = \text{succ}(\text{succ}(\text{double}(n)))$, as well as its proof by induction, also has to be “primed”.

In traditional mathematics, one just proclaims that \mathbb{N} and \mathbb{N}' are obviously “the same”, and can be substituted for each other whenever the need arises. This is usually unproblematic, but it sweeps a fair amount under the rug, widening the gap between informal mathematics and its precise description. In homotopy type theory, we can do better.

First observe that we have the following definable maps:

- $f := \text{rec}_{\mathbb{N}}(\mathbb{N}', 0', \lambda n. \text{succ}') : \mathbb{N} \rightarrow \mathbb{N}'$,
- $g := \text{rec}_{\mathbb{N}'}(\mathbb{N}, 0, \lambda n. \text{succ}) : \mathbb{N}' \rightarrow \mathbb{N}$.

Since the composition of g and f satisfies the same recurrences as the identity function on \mathbb{N} , Theorem 5.1.1 gives that $\prod_{(n:\mathbb{N})} g(f(n)) = n$, and the “primed” version of the same theorem gives $\prod_{(n:\mathbb{N}')} f(g(n)) = n$. Thus, f and g are quasi-inverses, so that $\mathbb{N} \simeq \mathbb{N}'$. We can now transfer functions on \mathbb{N} directly to functions on \mathbb{N}' (and vice versa) along this equivalence, e.g.

$$\text{double}' := \lambda n. f(\text{double}(g(n))).$$

It is an easy exercise to show that this version of double' is equal to the earlier one.

Of course, there is nothing surprising about this; such an isomorphism is exactly how a mathematician will envision “identifying” \mathbb{N} with \mathbb{N}' . However, the mechanism of “transfer” across

an isomorphism depends on the thing being transferred; it is not always as simple as pre- and post-composing a single function with f and g . Consider, for instance, a simple lemma such as

$$\prod_{n:\mathbb{N}'} \text{double}'(\text{succ}'(n)) = \text{succ}'(\text{succ}'(\text{double}'(n))).$$

Inserting the correct fs and gs is only a little easier than re-proving it by induction on $n : \mathbb{N}'$ directly.

Here is where the univalence axiom steps in: since $\mathbb{N} \simeq \mathbb{N}'$, we also have $\mathbb{N} =_{\mathcal{U}} \mathbb{N}'$, i.e. \mathbb{N} and \mathbb{N}' are *equal* as types. Now the induction principle for identity guarantees that any construction or proof relating to \mathbb{N} can automatically be transferred to \mathbb{N}' in the same way. We simply consider the type of the function or theorem as a type-indexed family of types $P : \mathcal{U} \rightarrow \mathcal{U}$, with the given object being an element of $P(\mathbb{N})$, and transport along the path $\mathbb{N} = \mathbb{N}'$. This involves considerably less overhead.

For simplicity, we have described this method in the case of two types \mathbb{N} and \mathbb{N}' with *identical*-looking definitions. However, a more common situation in practice is when the definitions are not literally identical, but nevertheless one induction principle implies the other. Consider, for instance, the type of lists from a one-element type, $\text{List}(\mathbf{1})$, which is generated by

- an element $\text{nil} : \text{List}(\mathbf{1})$, and
- a function $\text{cons} : \mathbf{1} \times \text{List}(\mathbf{1}) \rightarrow \text{List}(\mathbf{1})$.

This is not identical to the definition of \mathbb{N} , and it does not give rise to an identical induction principle. The induction principle of $\text{List}(\mathbf{1})$ says that for any $E : \text{List}(\mathbf{1}) \rightarrow \mathcal{U}$ together with recurrence data $e_{\text{nil}} : E(\text{nil})$ and $e_{\text{cons}} : \prod_{(u:\mathbf{1})} \prod_{(\ell:\text{List}(\mathbf{1}))} E(\ell) \rightarrow E(\text{cons}(u, \ell))$, there exists $f : \prod_{(\ell:\text{List}(\mathbf{1}))} E(\ell)$ such that $f(\text{nil}) \equiv e_{\text{nil}}$ and $f(\text{cons}(u, \ell)) \equiv e_{\text{cons}}(u, \ell, f(\ell))$. (We will see how to derive the induction principle of an inductive definition in §5.6.)

Now suppose we define $0'' := \text{nil} : \text{List}(\mathbf{1})$, and $\text{succ}'' : \text{List}(\mathbf{1}) \rightarrow \text{List}(\mathbf{1})$ by $\text{succ}''(\ell) := \text{cons}(\star, \ell)$. Then for any $E : \text{List}(\mathbf{1}) \rightarrow \mathcal{U}$ together with $e_0 : E(0'')$ and $e_s : \prod_{(\ell:\text{List}(\mathbf{1}))} E(\ell) \rightarrow E(\text{succ}''(\ell))$, we can define

$$\begin{aligned} e_{\text{nil}} &:= e_0 \\ e_{\text{cons}}(\star, \ell, x) &:= e_s(\ell, x). \end{aligned}$$

(In the definition of e_{cons} we use the induction principle of $\mathbf{1}$ to assume that u is \star .) Now we can apply the induction principle of $\text{List}(\mathbf{1})$, obtaining $f : \prod_{(\ell:\text{List}(\mathbf{1}))} E(\ell)$ such that

$$\begin{aligned} f(0'') &\equiv f(\text{nil}) \equiv e_{\text{nil}} \equiv e_0 \\ f(\text{succ}''(\ell)) &\equiv f(\text{cons}(\star, \ell)) \equiv e_{\text{cons}}(\star, \ell, f(\ell)) \equiv e_s(\ell, f(\ell)). \end{aligned}$$

Thus, $\text{List}(\mathbf{1})$ satisfies the same induction principle as \mathbb{N} , and hence (by the same arguments above) is equal to it.

Finally, these conclusions are not confined to the natural numbers: they apply to any inductive type. If we have an inductively defined type W , say, and some other type W' which satisfies the same induction principle as W , then it follows that $W \simeq W'$, and hence $W = W'$. We use the derived recursion principles for W and W' to construct maps $W \rightarrow W'$ and $W' \rightarrow W$, respectively, and then the induction principles for each to prove that both composites are equal to identities. For instance, in Chapter 1 we saw that the coproduct $A + B$ could also have been defined as $\sum_{(x:2)} \text{rec}_2(\mathcal{U}, A, B, x)$. The latter type satisfies the same induction principle as the former; hence they are canonically equivalent.

This is, of course, very similar to the familiar fact in category theory that if two objects have the same *universal property*, then they are equivalent. In §5.4 we will see that inductive types actually do have a universal property, so that this is a manifestation of that general principle.

5.3 W-types

Inductive types are very general, which is excellent for their usefulness and applicability, but makes them difficult to study as a whole. Fortunately, they can all be formally reduced to a few special cases. It is beyond the scope of this book to discuss this reduction — which is anyway irrelevant to the mathematician using type theory in practice — but we will take a little time to discuss the one of the basic special cases that we have not yet met. These are Martin-Löf’s *W-types*, also known as the types of *well-founded trees*. W-types are a generalization of such types as natural numbers, lists, and binary trees, which are sufficiently general to encapsulate the “recursion” aspect of *any* inductive type.

A particular W-type is specified by giving two parameters $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, in which case the resulting W-type is written $W_{(a:A)}B(a)$. The type A represents the type of *labels* for $W_{(a:A)}B(a)$, which function as constructors (however, we reserve that word for the actual functions which arise in inductive definitions). For instance, when defining natural numbers as a W-type, the type A would be the type **2** inhabited by the two elements 0_2 and 1_2 , since there are precisely two ways to obtain a natural number — either it will be zero or a successor of another natural number.

The type family $B : A \rightarrow \mathcal{U}$ is used to record the arity of labels: a label $a : A$ will take a family of inductive arguments, indexed over $B(a)$. We can therefore think of the “ $B(a)$ -many” arguments of a . These arguments are represented by a function $f : B(a) \rightarrow W_{(a:A)}B(a)$, with the understanding that for any $b : B(a)$, $f(b)$ is the “ b -th” argument to the label a . The W-type $W_{(a:A)}B(a)$ can thus be thought of as the type of well-founded trees, where nodes are labeled by elements of A and each node labeled by $a : A$ has $B(a)$ -many branches.

In the case of natural numbers, the label 0_2 has arity 0, since it constructs the constant zero; the label 1_2 has arity 1, since it constructs the successor of its argument. We can capture this by using simple elimination on **2** to define a function $\text{rec}_2(\mathcal{U}, \mathbf{0}, \mathbf{1})$ into a universe of types; this function returns the empty type **0** for 0_2 and the unit type **1** for 1_2 . We can thus define

$$\mathbf{N}^{\mathbf{w}} := W_{(b:\mathbf{2})}\text{rec}_2(\mathcal{U}, \mathbf{0}, \mathbf{1}, b)$$

where the superscript **w** serves to distinguish this version of natural numbers from the previously used one. Similarly, we can define the type of lists over A as a W-type with **1 + A** many labels: one nullary label for the empty list, plus one unary label for each $a : A$, corresponding to appending a to the head of a list:

$$\text{List}(A) := W_{(x:\mathbf{1+A})}\text{rec}_{\mathbf{1+A}}(\mathcal{U}, \mathbf{0}, \lambda a. \mathbf{1}, x).$$

In general, the W-type $W_{(x:A)}B(x)$ specified by $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$ is the inductive type generated by the following constructor:

- $\text{sup} : \prod_{(a:A)} (B(a) \rightarrow W_{(x:A)}B(x)) \rightarrow W_{(x:A)}B(x).$

The constructor sup (short for supremum) takes a label $a : A$ and a function $f : B(a) \rightarrow W_{(x:A)}B(x)$ representing the arguments to a , and constructs a new element of $W_{(x:A)}B(x)$. Using

our previous encoding of natural numbers as W-types, we can for instance define

$$0^W := \text{sup}(0_2, \lambda x. \text{rec}_0(N^W, x)).$$

Put differently, we use the label 0_2 to construct 0^W . Then, $\text{rec}_2(\mathcal{U}, \mathbf{0}, \mathbf{1}, 0_2)$ evaluates to $\mathbf{0}$, as it should since 0_2 is a nullary label. Thus, we need to construct a function $f : \mathbf{0} \rightarrow N^W$, which represents the (zero) arguments supplied to 0_2 . This is of course trivial, using simple elimination on $\mathbf{0}$ as shown. Similarly, we can define 1^W and a successor function succ^W

$$\begin{aligned} 1^W &:= \text{sup}(1_2, \lambda x. 0^W) \\ \text{succ}^W &:= \lambda n. \text{sup}(1_2, \lambda x. n). \end{aligned}$$

We have the following induction principle for W-types:

- When proving a statement $E : (W_{(x:A)} B(x)) \rightarrow \mathcal{U}$ about all elements of the W-type $W_{(x:A)} B(x)$, it suffices to prove it for $\text{sup}(a, f)$, assuming it holds for all $f(b)$ with $b : B(a)$. In other words, it suffices to give a proof

$$e : \prod_{(a:A)} \prod_{(f:B(a) \rightarrow W_{(x:A)} B(x))} \prod_{(g:\prod_{(b:B(a))} E(f(b)))} E(\text{sup}(a, f))$$

The variable g represents our inductive hypothesis, namely that all arguments of a satisfy E . To state this, we quantify over all elements of type $B(a)$, since each $b : B(a)$ corresponds to one argument $f(b)$ of a .

How would we define the function double on natural numbers encoded as a W-type? We would like to use the recursion principle of N^W with a codomain of N^W itself. We thus need to construct a suitable function

$$e : \prod_{(a:2)} \prod_{(f:B(a) \rightarrow N^W)} \prod_{(g:B(a) \rightarrow N^W)} N^W$$

which will represent the recurrence for the double function; for simplicity we denote the type family $\text{rec}_2(\mathcal{U}, \mathbf{0}, \mathbf{1})$ by B .

Clearly, e will be a function taking $a : 2$ as its first argument. The next step is to perform case analysis on a and proceed based on whether it is 0_2 or 1_2 . This suggests the following form

$$e := \lambda a. \text{rec}_2(C, e_0, e_1, a)$$

where

$$C := \prod_{(f:B(a) \rightarrow N^W)} \prod_{(g:B(a) \rightarrow N^W)} N^W.$$

If a is 0_2 , the type $B(a)$ becomes $\mathbf{0}$. Thus, given $f : \mathbf{0} \rightarrow N^W$ and $g : \mathbf{0} \rightarrow N^W$, we want to construct an element of N^W . Since the label 0_2 represents $\mathbf{0}$, it needs zero inductive arguments and the variables f and g are irrelevant. We return 0^W as a result:

$$e_0 := \lambda f. \lambda g. 0^W.$$

Analogously, if a is 1_2 , the type $B(a)$ becomes $\mathbf{1}$. Since the label 1_2 represents the successor operator, it needs one inductive argument — the predecessor — which is represented by the variable $f : \mathbf{1} \rightarrow N^W$. The value of the recursive call on the predecessor is represented by the

variable $g : \mathbf{1} \rightarrow \mathbf{N}^w$. Thus, taking this value (namely $g(\star)$) and applying the successor function twice thus yields the desired result:

$$e_1 := \lambda f. \lambda g. \text{succ}^w(\text{succ}^w(g(\star))).$$

Putting this together, we thus have

$$\text{double} := \text{rec}_{\mathbf{N}^w}(\mathbf{N}^w, e)$$

with e as defined above.

The associated computation rule for the function $\text{rec}_{W_{(x:A)}B(x)}(E, e) : \prod_{(w:W_{(x:A)}B(x))} E(w)$ is as follows.

- For any $a : A$ and $f : B(a) \rightarrow W_{(x:A)}B(x)$ we have

$$\text{rec}_{W_{(x:A)}B(x)}(E, e, \text{sup}(a, f)) \equiv e(a, f, (\lambda b. \text{rec}_{W_{(x:A)}B(x)}(E, e, f(b)))).$$

In other words, the function $\text{rec}_{W_{(x:A)}B(x)}(E, e)$ satisfies the recurrence e .

By the above computation rule, the function `double` behaves as expected:

$$\begin{aligned} \text{double}(0^w) &\equiv \text{rec}_{\mathbf{N}^w}(\mathbf{N}^w, e, \text{sup}(0_2, \lambda x. \text{rec}_0(\mathbf{N}^w, x))) \\ &\equiv e(0_2, (\lambda x. \text{rec}_0(\mathbf{N}^w, x)), (\lambda x. \text{double}(\text{rec}_0(\mathbf{N}^w, x)))) \\ &\equiv e_0((\lambda x. \text{rec}_0(\mathbf{N}^w, x)), (\lambda x. \text{double}(\text{rec}_0(\mathbf{N}^w, x)))) \\ &\equiv 0^w \end{aligned}$$

and

$$\begin{aligned} \text{double}(1^w) &\equiv \text{rec}_{\mathbf{N}^w}(\mathbf{N}^w, e, \text{sup}(1_2, \lambda x. 0^w)) \\ &\equiv e(1_2, (\lambda x. 0^w), (\lambda x. \text{double}(0^w))) \\ &\equiv e_1((\lambda x. 0^w), (\lambda x. \text{double}(0^w))) \\ &\equiv \text{succ}^w(\text{succ}^w((\lambda x. \text{double}(0^w))(\star))) \\ &\equiv \text{succ}^w(\text{succ}^w(0^w)) \end{aligned}$$

and so on.

Just as for natural numbers, we can prove a uniqueness theorem for W -types:

Theorem 5.3.1. *Let $g, h : \prod_{(w:W_{(x:A)}B(x))} E(w)$ be two functions which satisfy the recurrence*

$$e : \prod_{a,f} \left(\prod_{b:B(a)} E(f(b)) \right) \rightarrow E(\text{sup}(a, f)),$$

propositionally, i.e., such that

$$\begin{aligned} \prod_{a,f} g(\text{sup}(a, f)) &= e(a, f, \lambda b. g(f(b))), \\ \prod_{a,f} h(\text{sup}(a, f)) &= e(a, f, \lambda b. h(f(b))). \end{aligned}$$

Then g and h are equal.

5.4 Inductive types are initial algebras

As suggested earlier, inductive types also have a category-theoretic universal property. They are *homotopy-initial algebras*: initial objects (up to coherent homotopy) in a category of “algebras” determined by the specified constructors. As a simple example, consider the natural numbers. The appropriate sort of “algebra” here is a type equipped with the same structure that the constructors of \mathbb{N} give to it.

Definition 5.4.1. A \mathbb{N} -**algebra** is a type C with two elements $c_0 : C$, $c_s : C \rightarrow C$. The type of such algebras is

$$\mathbb{N}\text{Alg} := \sum_{C:\mathcal{U}} C \times (C \rightarrow C).$$

Definition 5.4.2. A \mathbb{N} -**homomorphism** between \mathbb{N} -algebras (C, c_0, c_s) and (D, d_0, d_s) is a function $h : C \rightarrow D$ such that $h(c_0) = d_0$ and $h(c_s(c)) = d_s(h(c))$ for all $c : C$. The type of such homomorphisms is

$$\mathbb{N}\text{Hom}((C, c_0, c_s), (D, d_0, d_s)) := \sum_{(h:C \rightarrow D)} (h(c_0) = d_0) \times \prod_{(c:C)} (h(c_s(c)) = d_s(h(c))).$$

We thus have a category of \mathbb{N} -algebras and \mathbb{N} -homomorphisms, and the claim is that \mathbb{N} is the initial object of this category. A category theorist will immediately recognize this as the definition of a *natural numbers object* in a category.

Of course, since our types behave like ∞ -groupoids, we actually have an $(\infty, 1)$ -category of \mathbb{N} -algebras, and we should ask \mathbb{N} to be initial in the appropriate $(\infty, 1)$ -categorical sense. Fortunately, we can formulate this without needing to define $(\infty, 1)$ -categories.

Definition 5.4.3. A \mathbb{N} -algebra I is called **homotopy-initial**, or **h-initial** for short, if for any other \mathbb{N} -algebra C , the type of \mathbb{N} -homomorphisms from I to C is contractible. Thus,

$$\text{isHinit}_{\mathbb{N}}(I) := \prod_{C:\mathbb{N}\text{Alg}} \text{isContr}(\mathbb{N}\text{Hom}(I, C)).$$

When they exist, h-initial algebras are unique — not just up to isomorphism, as usual in category theory, but up to equality, by the univalence axiom.

Theorem 5.4.4. *Any two h-initial \mathbb{N} -algebras are equal. Thus, the type of h-initial \mathbb{N} -algebras is a mere proposition.*

Proof. Suppose I and J are h-initial \mathbb{N} -algebras. Then $\mathbb{N}\text{Hom}(I, J)$ is contractible, hence inhabited by some \mathbb{N} -homomorphism $f : I \rightarrow J$, and likewise we have an \mathbb{N} -homomorphism $g : J \rightarrow I$. Now the composite $g \circ f$ is a \mathbb{N} -homomorphism from I to I , as is id_I ; but $\mathbb{N}\text{Hom}(I, I)$ is contractible, so $g \circ f = \text{id}_I$. Similarly, $f \circ g = \text{id}_J$. Hence $I \simeq J$, and so $I = J$. Since being contractible is a mere proposition and dependent products preserve mere propositions, it follows that being h-initial is itself a mere proposition. Thus any two proofs that I (or J) is h-initial are necessarily equal, which finishes the proof. \square

We now have the following theorem.

Theorem 5.4.5. *The \mathbb{N} -algebra $(\mathbb{N}, \mathbf{0}, \text{succ})$ is homotopy initial.*

Sketch of proof. Fix an arbitrary \mathbb{N} -algebra (C, c_0, c_s) . The recursion principle of \mathbb{N} yields a function $f : \mathbb{N} \rightarrow C$ defined by

$$\begin{aligned} f(0) &:= c_0 \\ f(\text{succ}(n)) &:= c_s(f(n)). \end{aligned}$$

These two equalities make f an \mathbb{N} -homomorphism, which we can take as the center of contraction for $\mathbb{N}\text{Hom}(\mathbb{N}, C)$. The uniqueness theorem (Theorem 5.1.1) then implies that any other \mathbb{N} -homomorphism is equal to f . \square

To place this in a more general context, it is useful to consider the notion of *algebra for an endofunctor*. Note that to make a type C into a \mathbb{N} -algebra is the same as to give a function $c : C + \mathbf{1} \rightarrow C$, and a function $f : C \rightarrow D$ is a \mathbb{N} -homomorphism just when $f \circ c \sim d \circ (f + \mathbf{1})$. In categorical language, this means the \mathbb{N} -algebras are the algebras for the endofunctor $F(X) := X + 1$ of the category of types.

For a more generic case, consider the \mathbb{W} -type associated to $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$. In this case we have an associated **polynomial functor**:

$$P(X) = \sum_{x:A} (B(x) \rightarrow X). \quad (5.4.6)$$

Actually, this assignment is functorial only up to homotopy, but this makes no difference in what follows. By definition, a **P -algebra** is then a type C equipped with a function $s_C : PC \rightarrow C$. By the universal property of Σ -types, this is equivalent to giving a function $\prod_{(a:A)} (B(a) \rightarrow C) \rightarrow C$. We will also call such objects **\mathbb{W} -algebras** for A and B , and we write

$$\mathbb{W}\text{Alg}(A, B) := \sum_{(C:\mathcal{U})} \prod_{(a:A)} (B(a) \rightarrow C) \rightarrow C.$$

Similarly, for P -algebras (C, s_C) and (D, s_D) , a **homomorphism** between them $(f, s_f) : (C, s_C) \rightarrow (D, s_D)$ consists of a function $f : C \rightarrow D$ and a homotopy between maps $PC \rightarrow PD$

$$s_f : f \circ s_C = s_D \circ Pf,$$

where $Pf : PC \rightarrow PD$ is the result of the easily-definable action of P on $f : C \rightarrow D$. Such an algebra homomorphism can be represented suggestively in the form:

$$\begin{array}{ccc} PC & \xrightarrow{Pf} & PD \\ s_C \downarrow & s_f & \downarrow s_D \\ C & \xrightarrow{f} & D \end{array}$$

In terms of elements, f is a P -homomorphism (or **\mathbb{W} -homomorphism**) if

$$f(s_C(a, h)) = s_D(a, f \circ h).$$

We have the type of \mathbb{W} -homomorphisms:

$$\mathbb{W}\text{Hom}_{A,B}((C, s_C), (D, s_D)) := \sum_{(f:C \rightarrow D)} \prod_{(a:A)} \prod_{(h:B(a) \rightarrow C)} f(s_C(a, h)) = s_D(a, f \circ h)$$

Finally, a P -algebra (C, s_C) is said to be **homotopy-initial** if for every P -algebra (D, s_D) , the type of all algebra homomorphisms $(C, s_C) \rightarrow (D, s_D)$ is contractible. That is,

$$\text{isHinit}_W(A, B, I) := \prod_{C: W\text{Alg}(A, B)} \text{isContr}(\text{WHom}_{A, B}(I, C)).$$

Now the analogous theorem to Theorem 5.4.5 is:

Theorem 5.4.7. *For any type $A : \mathcal{U}$ and type family $B : A \rightarrow \mathcal{U}$, the W -algebra $(W_{(x:A)} B(x), \sup)$ is h -initial.*

Sketch of proof. Suppose we have $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, and consider the associated polynomial functor $P(X) := \sum_{(x:A)} (B(x) \rightarrow X)$. Let $W := W_{(x:A)} B(x)$. Then using the W -introduction rule from §5.3, we have a structure map $s_W := \sup : PW \rightarrow W$. We want to show that the algebra (W, s_W) is h -initial. So, let us consider another algebra (C, s_C) and show that the type $T := \text{WHom}_{A, B}((W, s_W), (C, s_C))$ of W -homomorphisms from (W, s_W) to (C, s_C) is contractible. To do so, observe that the W -elimination rule and the W -computation rule allow us to define a W -homomorphism $(f, s_f) : (W, s_W) \rightarrow (C, s_C)$, thus showing that T is inhabited. It is furthermore necessary to show that for every W -homomorphism $(g, s_g) : (W, s_W) \rightarrow (C, s_C)$, there is an identity proof

$$p : (f, s_f) = (g, s_g). \quad (5.4.8)$$

This uses the fact that, in general, a type of the form $(f, s_f) = (g, s_g)$ is equivalent to the type of what we call **algebra 2-cells** from f to g , whose canonical elements are pairs of the form (e, s_e) , where $e : f = g$ and s_e is a higher identity proof between the identity proofs represented by the following pasting diagrams:

In light of this fact, to prove that there exists an element as in (5.4.8), it is sufficient to show that there is an algebra 2-cell (e, s_e) from f to g . The identity proof $e : f = g$ is now constructed by function extensionality and W -elimination so as to guarantee the existence of the required identity proof s_e . \square

5.5 Homotopy-inductive types

In §5.3 we showed how to encode natural numbers as W -types, with

$$\begin{aligned} \mathbf{N}^W &:= W_{(b:2)} \text{rec}_2(\mathcal{U}, \mathbf{0}, \mathbf{1}, b), \\ \mathbf{0}^W &:= \sup(0_2, (\lambda x. \text{rec}_0(\mathbf{N}^W, x))), \\ \text{succ}^W &:= \lambda n. \sup(1_2, (\lambda x. n)). \end{aligned}$$

We also showed how one can define a double function on \mathbf{N}^W using the recursion principle. When it comes to the induction principle, however, this encoding is no longer satisfactory: given $E :$

$\mathbf{N}^{\mathbf{W}} \rightarrow \mathcal{U}$ and recurrences $e_z : E(0^{\mathbf{W}})$ and $e_s : \prod_{(n:\mathbf{N}^{\mathbf{W}})} E(n) \rightarrow E(\text{succ}^{\mathbf{W}}(n))$, we can only construct a dependent function $r(E, e_z, e_s) : \prod_{(n:\mathbf{N}^{\mathbf{W}})} E(n)$ satisfying the given recurrences *propositionally*, i.e. up to a path. This means that the computation rules for natural numbers, which give judgmental equalities, cannot be derived from the rules for W-types in any obvious way.

This problem goes away if instead of the conventional inductive types we consider *homotopy-inductive types*, where all computation rules are stated up to a path, i.e. the symbol \equiv is replaced by $=$. For instance, the computation rule for the homotopy version of W-types \mathbf{W}^h becomes:

- For any $a : A$ and $f : B(a) \rightarrow \mathbf{W}_{(x:A)}^h B(x)$ we have

$$\text{rec}_{\mathbf{W}_{(x:A)}^h B(x)}(E, e, \sup(a, f)) = e\left(a, f, (\lambda b. \text{rec}_{\mathbf{W}_{(x:A)}^h B(x)}(E, f(b)))\right)$$

Homotopy-inductive types have an obvious disadvantage when it comes to computational properties — the behavior of any function constructed using the induction principle can now only be characterized propositionally. But numerous other considerations drive us to consider homotopy-inductive types as well. For instance, while we showed in §5.4 that inductive types are homotopy-initial algebras, not every homotopy-initial algebra is an inductive type (i.e. satisfies the corresponding induction principle) — but every homotopy-initial algebra *is* a homotopy-inductive type. Similarly, we might want to apply the uniqueness argument from §5.2 when one (or both) of the types involved is only a homotopy-inductive type — for instance, to show that the W-type encoding of \mathbb{N} is equivalent to the usual \mathbb{N} .

Additionally, the notion of a homotopy-inductive type is now internal to the type theory. For example, this means we can form a type of all natural numbers objects and make assertions about it. In the case of W-types, we can characterize a homotopy W-type $\mathbf{W}_{(x:A)} B(x)$ as any type endowed with a supremum function and an induction principle satisfying the appropriate (propositional) computation rule:

$$\begin{aligned} \mathbf{W}_d(A, B) := & \sum_{(W:\mathcal{U})} \sum_{(\sup:\prod_{(a:A)}(B(a) \rightarrow W) \rightarrow W)} \prod_{(E:W \rightarrow \mathcal{U})} \\ & \prod_{(e:\prod_{(a,f)}(\prod_{(b:B(a))} E(f(b)) \rightarrow E(\sup(a,f))))} \sum_{(\text{ind}:\prod_{(w:W)} E(w))} \prod_{(a,f)} \\ & \text{ind}(\sup(a,f)) = e(a, \lambda b. \text{ind}(f(b))). \end{aligned}$$

In Chapter 6 we will see some other reasons why propositional computation rules are worth considering.

In this section, we will state some basic facts about homotopy-inductive types. We omit most of the proofs, which are somewhat technical.

Theorem 5.5.1. *For any $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, the type $\mathbf{W}_d(A, B)$ is a mere proposition.*

It turns out that there is an equivalent characterization of W-types using a recursion principle, plus certain *uniqueness* and *coherence* laws. First we give the recursion principle:

- When constructing a function from the W-type $\mathbf{W}_{(x:A)}^h B(x)$ into the type C , it suffices to give its value for $\sup(a, f)$, assuming we are given the values of all $f(b)$ with $b : B(a)$. In other words, it suffices to construct a function

$$c : \prod_{a:A} (B(a) \rightarrow C) \rightarrow C.$$

The associated computation rule for $\text{rec}_{W_{(x:A)}^h B(x)}(C, c) : (W_{(x:A)} B(x)) \rightarrow C$ is as follows:

- For any $a : A$ and $f : B(a) \rightarrow W_{(x:A)}^h B(x)$ we have a witness $\beta(C, c, a, f)$ for equality

$$\text{rec}_{W_{(x:A)}^h B(x)}(C, c, \sup(a, f)) = c(a, \lambda b. \text{rec}_{W_{(x:A)}^h B(x)}(C, c, f(b))).$$

Furthermore, we assert the following uniqueness principle, saying that any two functions defined by the same recurrence are equal:

- Let $C : \mathcal{U}$ and $c : \prod_{(a:A)}(B(a) \rightarrow C) \rightarrow C$ be given. Let $g, h : (W_{(x:A)}^h B(x)) \rightarrow C$ be two functions which satisfy the recurrence c up to propositional equality, i.e., such that we have

$$\begin{aligned} \beta_g &: \prod_{a,f} g(\sup(a, f)) = c(a, \lambda b. g(f(b))), \\ \beta_h &: \prod_{a,f} h(\sup(a, f)) = c(a, \lambda b. h(f(b))). \end{aligned}$$

Then g and h are equal, i.e. there is $\alpha(C, c, f, g, \beta_g, \beta_h)$ of type $g = h$.

Recall that when we have an induction principle rather than only a recursion principle, this propositional uniqueness principle is derivable (Theorem 5.3.1). But with only recursion, the uniqueness principle is no longer derivable — and in fact, the statement is not even true (exercise). Hence, we postulate it as an axiom. We also postulate the following coherence law, which tells us how the proof of uniqueness behaves on canonical elements:

- For any $a : A$ and $f : B(a) \rightarrow C$, the following diagram commutes propositionally:

$$\begin{array}{ccc} g(\sup(x, f)) & \xrightarrow{\beta_g} & c(a, \lambda b. g(f(b))) \\ \alpha(\sup(x, f)) \downarrow & & \downarrow c(a, -)(\text{funext}(\lambda b. \alpha(f(b)))) \\ h(\sup(x, f)) & \xrightarrow{\beta_h} & c(a, \lambda b. h(f(b))) \end{array}$$

where α abbreviates the path $\alpha(C, c, f, g, \beta_g, \beta_h) : g = h$.

Putting all of this data together yields another characterization of $W_{(x:A)} B(x)$, as a type with a supremum function, satisfying simple elimination, computation, uniqueness, and coherence rules:

$$\begin{aligned} W_s(A, B) &:= \sum_{(W:\mathcal{U})} \sum_{(\sup:\prod_{(a:A)}(B(a) \rightarrow W) \rightarrow W)} \prod_{(C:\mathcal{U})} \prod_{(c:\prod_{(a:A)}(B(a) \rightarrow C) \rightarrow C)} \\ &\quad \sum_{(\text{rec}:W \rightarrow C)} \sum_{(\beta:\prod_{(a,f)} \text{rec}(\sup(a,f)) = c(a, \lambda b. \text{rec}(f(b))))} \prod_{(g:W \rightarrow C)} \prod_{(h:W \rightarrow C)} \prod_{(\beta_g:\prod_{(a,f)} g(\sup(a,f)) = c(a, \lambda b. g(f(b))))} \\ &\quad \prod_{(\beta_h:\prod_{(a,f)} h(\sup(a,f)) = c(a, \lambda b. h(f(b))))} \sum_{(\alpha:\prod_{(w:W)} g(w) = h(w))} \\ &\quad \alpha(\sup(x, f)) \cdot \beta_h = \beta_g \cdot c(a, -)(\text{funext} \lambda b. \alpha(f(b))) \end{aligned}$$

Theorem 5.5.2. *For any $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, the type $W_s(A, B)$ is a mere proposition.*

Finally, we have a third, very concise characterization of $W_{(x:A)} B(x)$ as an h-initial W -algebra:

$$W_h(A, B) := \sum_{I:W\text{Alg}(A, B)} \text{isHinit}_W(A, B, I).$$

Theorem 5.5.3. *For any $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, the type $\mathsf{W}_h(A, B)$ is a mere proposition.*

It turns out all three characterizations of W -types are in fact equivalent:

Lemma 5.5.4. *For any $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, we have*

$$\mathsf{W}_d(A, B) \simeq \mathsf{W}_s(A, B) \simeq \mathsf{W}_h(A, B)$$

Indeed, we have the following theorem, which is an improvement over Theorem 5.4.7:

Theorem 5.5.5. *The types satisfying the formation, introduction, elimination, and propositional computation rules for W -types are precisely the homotopy-initial W -algebras.*

Sketch of proof. Inspecting the proof of Theorem 5.4.7, we see that only the *propositional* computation rule was required to establish the h-initiality of $\mathsf{W}_{(x:A)}B(x)$. For the converse implication, let us assume that the polynomial functor associated to $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, has an h-initial algebra (W, s_W) ; we show that W satisfies the propositional rules of W -types. The W -introduction rule is simple; namely, for $a : A$ and $t : B(a) \rightarrow W$, we define $\text{sup}(a, t) : W$ to be the result of applying the structure map $s_W : PW \rightarrow W$ to $(a, t) : PW$. For the W -elimination rule, let us assume its premisses and in particular that $C' : W \rightarrow \mathcal{U}$. Using the other premisses, one shows that the type $C := \sum_{(w:W)} C'(w)$ can be equipped with a structure map $s_C : PC \rightarrow C$. By the h-initiality of W , we obtain an algebra homomorphism $(f, s_f) : (W, s_W) \rightarrow (C, s_C)$. Furthermore, the first projection $\text{pr}_1 : C \rightarrow W$ can be equipped with the structure of a homomorphism, so that we obtain a diagram of the form

$$\begin{array}{ccccc} PW & \xrightarrow{Pf} & PC & \xrightarrow{P\text{pr}_1} & PW \\ s_W \downarrow & & s_C \downarrow & & s_W \downarrow \\ W & \xrightarrow{f} & C & \xrightarrow{\text{pr}_1} & W. \end{array}$$

But the identity function $1_W : W \rightarrow W$ has a canonical structure of an algebra homomorphism and so, by the contractibility of the type of homomorphisms from (W, s_W) to itself, there must be an identity proof between the composite of (f, s_f) with $(\text{pr}_1, s_{\text{pr}_1})$ and $(1_W, s_{1_W})$. This implies, in particular, that there is an identity proof $p : \text{pr}_1 \circ f = 1_W$.

Since $(\text{pr}_2 \circ f)w : C((\text{pr}_1 \circ f)w)$, we can define

$$\text{rec}(w, c) := p_*((\text{pr}_2 \circ f)w) : C(w)$$

where the transport p_* is with respect to the family

$$\lambda u. C \circ u : (W \rightarrow W) \rightarrow W \rightarrow \mathcal{U}.$$

The verification of the propositional W -computation rule is a calculation, involving the naturality properties of operations of the form p_* . \square

Finally, as desired, we can encode homotopy-natural-numbers as homotopy- W -types:

Theorem 5.5.6. *The rules for natural numbers with propositional computation rules can be derived from the rules for W -types with propositional computation rules.*

5.6 The general syntax of inductive definitions

So far, we have been discussing only particular inductive types: **0**, **1**, **2**, \mathbb{N} , coproducts, products, Σ -types, W -types, etc. However, an important aspect of type theory is the ability to define *new* inductive types, rather than being restricted only to some particular fixed list of them. In order to be able to do this, however, we need to know what sorts of “inductive definitions” are valid or reasonable.

To see that not everything which “looks like an inductive definition” makes sense, consider the following “constructor” of a type C :

- $g : (C \rightarrow \mathbb{N}) \rightarrow C$.

The recursion principle for such a type C ought to say that given a type P , in order to construct a function $f : C \rightarrow P$, it suffices to consider the case when the input $c : C$ is of the form $g(\alpha)$ for some $\alpha : C \rightarrow \mathbb{N}$. Moreover, we would expect to be able to use the “recursive data” of f applied to α in some way. However, it is not at all clear how to “apply f to α ”, since both are functions with domain C .

We could write down a “recursion principle” for C by just supposing (unjustifiably) that there is some way to apply f to α and obtain a function $P \rightarrow \mathbb{N}$. Then the input to the recursion rule would ask for a type P together with a function

$$h : (C \rightarrow \mathbb{N}) \rightarrow (P \rightarrow \mathbb{N}) \rightarrow P \tag{5.6.1}$$

where the two arguments of h are α and “the result of applying f to α ”. However, what would the computation rule for the resulting function $f : C \rightarrow P$ be? Looking at other computation rules, we would expect something like “ $f(g(\alpha)) \equiv h(\alpha, f(\alpha))$ ” for $\alpha : C \rightarrow \mathbb{N}$, but as we have seen, “ $f(\alpha)$ ” does not make sense. The induction principle of C is even more problematic; it’s not even clear how to write down the hypotheses.

On the other hand, we could write down a different “recursion principle” for C by ignoring the “recursive” presence of C in the domain of α , considering it as merely an indexing type for a family of natural numbers. In this case the input would ask for a type P together with a function

$$h : (C \rightarrow \mathbb{N}) \rightarrow P,$$

so the type of the recursion principle would be $\text{rec}_C : \prod_{(P:\mathcal{U})} ((C \rightarrow \mathbb{N}) \rightarrow P) \rightarrow C \rightarrow P$, and similarly for the induction principle. Now it is possible to write down a computation rule, namely $\text{rec}_C(P, h, g(\alpha)) \equiv h(\alpha)$. However, the existence of a type C with this recursor and computation rule turns out to be inconsistent. See Exercises 5.7 to 5.10 for proofs of this and other variations.

This example suggests one restriction on inductive definitions: the domains of all the constructors must be *covariant functors* of the type being defined, so that we can “apply f to them” to get the result of the “recursive call”. In other words, if we replace all occurrences of the type being defined with a variable $X : \mathcal{U}$, then each domain of a constructor must be an expression that can be made into a covariant functor of X . This is the case for all the examples we have considered so far. For instance, with the constructor $\text{inl} : A \rightarrow A + B$, the relevant functor is constant at A (i.e. $X \mapsto A$), while for the constructor $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$, the functor is the identity functor ($X \mapsto X$).

However, this necessary condition is also not sufficient. Covariance prevents the inductive type from occurring on the left of a single function type, as in the argument $C \rightarrow \mathbb{N}$ of the

“constructor” g considered above, since this yields a contravariant functor rather than a covariant one. However, since the composite of two contravariant functors is covariant, *double* function types such as $((X \rightarrow \mathbb{N}) \rightarrow \mathbb{N})$ are once again covariant. This enables us to reproduce Cantorian-style paradoxes.

For instance, consider an “inductive type” D with the following constructor:

- $k : ((D \rightarrow \text{Prop}) \rightarrow \text{Prop}) \rightarrow D$.

Assuming such a type exists, we define functions

$$\begin{aligned} r &: D \rightarrow (D \rightarrow \text{Prop}) \rightarrow \text{Prop}, \\ f &: (D \rightarrow \text{Prop}) \rightarrow D, \\ p &: (D \rightarrow \text{Prop}) \rightarrow (D \rightarrow \text{Prop}) \rightarrow \text{Prop}, \end{aligned}$$

by

$$\begin{aligned} r(k(\theta)) &:= \theta, \\ f(\delta) &:= k(\lambda x. (x = \delta)), \\ p(\delta) &:= \lambda x. \delta(f(x)). \end{aligned}$$

Here r is defined by the recursion principle of D , while f and p are defined explicitly. Then for any $\delta : D \rightarrow \text{Prop}$, we have $r(f(\delta)) = \lambda x. (x = \delta)$.

In particular, therefore, if $f(\delta) = f(\delta')$, then we have a path $s : (\lambda x. (x = \delta)) = (\lambda x. (x = \delta'))$. Thus, $\text{happly}(s, \delta) : (\delta = \delta) = (\delta = \delta')$, and so in particular $\delta = \delta'$ holds. Hence, f is “injective” (although *a priori* D may not be a set). This already sounds suspicious — we have an “injection” of the “power set” of D into D — and with a little more work we can massage it into a contradiction.

Suppose given $\theta : (D \rightarrow \text{Prop}) \rightarrow \text{Prop}$, and define $\delta : D \rightarrow \text{Prop}$ by

$$\delta(d) := \exists(\gamma : D \rightarrow \text{Prop}). (f(\gamma) = d) \times \theta(\gamma). \quad (5.6.2)$$

We claim that $p(\delta) = \theta$. By function extensionality, it suffices to show $p(\delta)(\gamma) =_{\text{Prop}} \theta(\gamma)$ for any $\gamma : D \rightarrow \text{Prop}$. And by univalence, for this it suffices to show that each implies the other. Now by definition of p , we have

$$\begin{aligned} p(\delta)(\gamma) &\equiv \delta(f(\gamma)) \\ &\equiv \exists(\gamma' : D \rightarrow \text{Prop}). (f(\gamma') = f(\gamma)) \times \theta(\gamma'). \end{aligned}$$

Clearly this holds if $\theta(\gamma)$, since we may take $\gamma' := \gamma$. On the other hand, if we have γ' with $f(\gamma') = f(\gamma)$ and $\theta(\gamma')$, then $\gamma' = \gamma$ since f is injective, hence also $\theta(\gamma)$.

This completes the proof that $p(\delta) = \theta$. Thus, every element $\theta : (D \rightarrow \text{Prop}) \rightarrow \text{Prop}$ is the image under p of some element $\delta : D \rightarrow \text{Prop}$. However, if we define θ by a classic diagonalization:

$$\theta(\gamma) := \neg p(\gamma)(\gamma) \quad \text{for all } \gamma : D \rightarrow \text{Prop}$$

then from $\theta = p(\delta)$ we deduce $p(\delta)(\delta) = \neg p(\delta)(\delta)$. This is a contradiction: no proposition can be equivalent to its negation. (Supposing $P \Leftrightarrow \neg P$, if P , then $\neg P$, and so $\mathbf{0}$; hence $\neg P$, but then P , and so $\mathbf{0}$.)

Remark 5.6.3. There is a question of universe size to be addressed. In general, an inductive type must live in a universe that already contains all the types going into its definition. Thus if in the definition of D , the ambiguous notation Prop means $\text{Prop}_{\mathcal{U}}$, then we do not have $D : \mathcal{U}$ but only $D : \mathcal{U}'$ for some larger universe \mathcal{U}' with $\mathcal{U} : \mathcal{U}'$. In a predicative theory, therefore, the right-hand side of (5.6.2) lives in $\text{Prop}_{\mathcal{U}'}$, not $\text{Prop}_{\mathcal{U}}$. So this contradiction does require the propositional resizing axiom mentioned in §3.5.

This counterexample suggests that we should ban an inductive type from ever appearing on the left of an arrow in the domain of its constructors, even if that appearance is nested in other arrows so as to eventually become covariant. (Similarly, we also forbid it from appearing in the domain of a dependent function type.) This restriction is called **strict positivity** (ordinary “positivity” being essentially covariance), and it turns out to suffice.

In conclusion, therefore, a valid inductive definition of a type W consists of a list of *constructors*. Each constructor is assigned a type that is a function type taking some number (possibly zero) of inputs (possibly dependent on one another) and returning an element of W . Finally, we allow W itself to occur in the input types of its constructors, but only strictly positively. This essentially means that each argument of a constructor is either a type not involving W , or some iterated function type with codomain W . For instance, the following is a valid constructor type:

$$c : (A \rightarrow W) \rightarrow (B \rightarrow C \rightarrow W) \rightarrow D \rightarrow W \rightarrow W. \quad (5.6.4)$$

All of these function types can also be dependent functions (Π -types).¹

Note we require that an inductive definition is given by a *finite* list of constructors. This is simply because we have to write it down on the page. If we want an inductive type which behaves as if it has an infinite number of constructors, we can simply parametrize one constructor by some infinite type. For instance, a constructor such as $\mathbb{N} \rightarrow W \rightarrow W$ can be thought of as equivalent to countably many constructors of the form $W \rightarrow W$. (Of course, the infinity is now *internal* to the type theory, but this is as it should be for any foundational system.) Similarly, if we want a constructor that takes “infinitely many arguments”, we can allow it to take a family of arguments parametrized by some infinite type, such as $(\mathbb{N} \rightarrow W) \rightarrow W$ which takes an infinite sequence of elements of W .

Now, once we have such an inductive definition, what can we do with it? Firstly, there is a **recursion principle** stating that in order to define a function $f : W \rightarrow P$, it suffices to consider the case when the input $w : W$ arises from one of the constructors, allowing ourselves to recursively call f on the inputs to that constructor. For the example constructor (5.6.4), we would require P to be equipped with a function of type

$$d : (A \rightarrow W) \rightarrow (A \rightarrow P) \rightarrow (B \rightarrow C \rightarrow W) \rightarrow (B \rightarrow C \rightarrow P) \rightarrow D \rightarrow W \rightarrow P \rightarrow P. \quad (5.6.5)$$

Under these hypotheses, the recursion principle yields $f : W \rightarrow P$, which moreover “preserves the constructor data” in the evident way — this is the computation rule, where we use covariance of the inputs. For instance, in the example (5.6.4), the computation rule says that for any $\alpha : A \rightarrow W$, $\beta : B \rightarrow C \rightarrow W$, $\delta : D$, and $\omega : W$, we have

$$f(c(\alpha, \beta, \delta, \omega)) \equiv d(\alpha, f \circ \alpha, \beta, f \circ \beta, \delta, \omega, f(\omega)). \quad (5.6.6)$$

¹In the language of §5.4, the condition of strict positivity ensures that the relevant endofunctor is polynomial. It is well-known in category theory that not *all* endofunctors can have initial algebras; restricting to polynomial functors ensures consistency. One can consider various relaxations of this condition, but in this book we will restrict ourselves to strict positivity as defined here.

The **induction principle** for a general inductive type W is only a little more complicated. Of course, we start with a type family $P : W \rightarrow \mathcal{U}$, which we require to be equipped with constructor data “lying over” the constructor data of W . That means the “recursive call” arguments such as $A \rightarrow P$ above must be replaced by dependent functions with types such as $\prod_{(a:A)} P(\alpha(a))$. In the full example of (5.6.4), the corresponding hypothesis for the induction principle would require

$$\begin{aligned} d : \prod_{\alpha:A \rightarrow W} \left(\prod_{a:A} P(\alpha(a)) \right) \rightarrow \prod_{\beta:B \rightarrow C \rightarrow W} \left(\prod_{(b:B)} \prod_{(c:C)} P(\beta(b,c)) \right) \rightarrow \\ \prod_{(\delta:D)} \prod_{(\omega:W)} P(\omega) \rightarrow P(c(\alpha, \beta, \delta, \omega)). \end{aligned} \quad (5.6.7)$$

The corresponding computation rule looks identical to (5.6.6). Of course, the recursion principle is the special case of the induction principle where P is a constant family. As we have mentioned before, the induction principle is also called the **eliminator**, and the recursion principle the **non-dependent eliminator**.

As discussed in §1.10, we also allow ourselves to invoke the induction and recursion principles implicitly, writing a definitional equation with \equiv for each expression that would be the hypotheses of the induction principle. This is called giving a definition by (dependent) **pattern matching**. In our running example, this means we could define $f : \prod_{(w:W)} P(w)$ by

$$f(c(\alpha, \beta, \delta, \omega)) \equiv \dots$$

where $\alpha : A \rightarrow W$ and $\beta : B \rightarrow C \rightarrow W$ and $\delta : D$ and $\omega : W$ are variables that are bound in the right-hand side. Moreover, the right-hand side may involve recursive calls to f of the form $f(\alpha(a))$, $f(\beta(b,c))$, and $f(\omega)$. When this definition is repackaged in terms of the induction principle, we replace such recursive calls by $\bar{\alpha}(a)$, $\bar{\beta}(b,c)$, and $\bar{\omega}$, respectively, for new variables

$$\begin{aligned} \bar{\alpha} &: \prod_{a:A} P(\alpha(a)) \\ \bar{\beta} &: \prod_{(b:B)} \prod_{(c:C)} P(\beta(b,c)) \\ \bar{\omega} &: P(\omega). \end{aligned}$$

Then we could write

$$f \equiv \text{ind}_W(P, \lambda\alpha. \lambda\bar{\alpha}. \lambda\beta. \lambda\bar{\beta}. \lambda\delta. \lambda\bar{\omega}. \dots)$$

where the second argument to ind_W has the type of (5.6.7).

We will not attempt to give a formal presentation of the grammar of a valid inductive definition and its resulting induction and recursion principles and pattern matching rules. This is possible to do (indeed, it is necessary to do if implementing a computer proof assistant), but provides no additional insight. With practice, one learns to automatically deduce the induction and recursion principles for any inductive definition, and to use them without having to think twice.

5.7 Generalizations of inductive types

The notion of inductive type has been studied in type theory for many years, and admits many, many generalizations: inductive type families, mutual inductive types, inductive-inductive types, inductive-recursive types, etc. In this section we give an overview of some of these, a few of

which will be used later in the book. (In Chapter 6 we will study in more depth a very different generalization of inductive types, which is particular to *homotopy type theory*.)

Most of these generalizations involve allowing ourselves to define more than one type by induction at the same time. One very simple example of this, which we have already seen, is the coproduct $A + B$. It would be tedious indeed if we had to write down separate inductive definitions for $\mathbb{N} + \mathbb{N}$, for $\mathbb{N} + 2$, for $2 + 2$, and so on every time we wanted to consider the coproduct of two types. Instead, we make one definition in which A and B are variables standing for types; in type theory they are called **parameters**. Thus technically speaking, what results from the definition is not a single type, but a family of types $+ : \mathcal{U} \rightarrow \mathcal{U} \rightarrow \mathcal{U}$, taking two types as input and producing their coproduct. Similarly, the type $\text{List}(A)$ of lists is a family $\text{List}(-) : \mathcal{U} \rightarrow \mathcal{U}$ in which the type A is a parameter.

In mathematics, this sort of thing is so obvious as to not be worth mentioning, but we bring it up in order to contrast it with the next example. Note that each type $A + B$ is *independently* defined inductively, as is each type $\text{List}(A)$. By contrast, we might also consider defining a whole type family $B : A \rightarrow \mathcal{U}$ by induction *together*. The difference is that now the constructors may change the index $a : A$, and as a consequence we cannot say that the individual types $B(a)$ are inductively defined, only that the entire family is inductively defined.

The standard example is the type of *lists of specified length*, traditionally called **vectors**. We fix a parameter type A , and define a type family $\text{Vec}_n(A)$, for $n : \mathbb{N}$, generated by the following constructors:

- a vector $\text{nil} : \text{Vec}_0(A)$ of length zero,
- a function $\text{cons} : \prod_{(n:\mathbb{N})} A \rightarrow \text{Vec}_n(A) \rightarrow \text{Vec}_{\text{succ}(n)}(A)$.

In contrast to lists, vectors (with elements from a fixed type A) form a family of types indexed by their length. While A is a parameter, we say that $n : \mathbb{N}$ is an **index** of the inductive family. An individual type such as $\text{Vec}_3(A)$ is not inductively defined: the constructors which build elements of $\text{Vec}_3(A)$ take input from a different type in the family, such as $\text{cons} : A \rightarrow \text{Vec}_2(A) \rightarrow \text{Vec}_3(A)$.

In particular, the induction principle must refer to the entire type family as well; thus the hypotheses and the conclusion must quantify over the indices appropriately. In the case of vectors, the induction principle states that given a type family $C : \prod_{(n:\mathbb{N})} \text{Vec}_n(A) \rightarrow \mathcal{U}$, together with

- an element $c_{\text{nil}} : C(0, \text{nil})$, and
- a function $c_{\text{cons}} : \prod_{(n:\mathbb{N})} \prod_{(a:A)} \prod_{(\ell:\text{Vec}_n(A))} C(n, \ell) \rightarrow C(\text{succ}(n), \text{cons}(a, \ell))$

there exists a function $f : \prod_{(n:\mathbb{N})} \prod_{(\ell:\text{Vec}_n(A))} C(n, \ell)$ such that

$$\begin{aligned} f(0, \text{nil}) &\equiv c_{\text{nil}} \\ f(\text{succ}(n), \text{cons}(a, \ell)) &\equiv c_{\text{cons}}(n, a, \ell, f(\ell)). \end{aligned}$$

One use of inductive families is to define *predicates* inductively. For instance, we might define the predicate $\text{iseven} : \mathbb{N} \rightarrow \mathcal{U}$ as an inductive family indexed by \mathbb{N} , with the following constructors:

- an element $\text{even}_0 : \text{iseven}(0)$,
- a function $\text{even}_{ss} : \prod_{(n:\mathbb{N})} \text{iseven}(n) \rightarrow \text{iseven}(\text{succ}(\text{succ}(n)))$.

In other words, we stipulate that 0 is even, and that if n is even then so is $\text{succ}(\text{succ}(n))$. These constructors “obviously” give no way to construct an element of, say, $\text{iseven}(1)$, and since iseven is supposed to be freely generated by these constructors, there must be no such element. (Actually proving that $\neg \text{iseven}(1)$ is not entirely trivial, however). The induction principle for iseven says that to prove something about all even natural numbers, it suffices to prove it for 0 and verify that it is preserved by adding two.

Inductively defined predicates are much used in computer formalization of mathematics and software verification. But we will not have much use for them, with a couple of exceptions in **????**.

Another important special case is when the indexing type of an inductive family is finite. In this case, we can equivalently express the inductive definition as a finite collection of types defined by *mutual induction*. For instance, we might define the types even and odd of even and odd natural numbers by mutual induction, where even is generated by constructors

- $0 : \text{even}$ and
- $\text{esucc} : \text{odd} \rightarrow \text{even}$,

while odd is generated by the one constructor

- $\text{osucc} : \text{even} \rightarrow \text{odd}$.

Note that even and odd are simple types (not type families), but their constructors can refer to each other. If we expressed this definition as an inductive type family $\text{paritynat} : 2 \rightarrow \mathcal{U}$, with $\text{paritynat}(0_2)$ and $\text{paritynat}(1_2)$ representing even and odd respectively, it would instead have constructors:

- $0 : \text{paritynat}(0_2)$,
- $\text{esucc} : \text{paritynat}(1_2) \rightarrow \text{paritynat}(0_2)$,
- $\text{osucc} : \text{paritynat}(0_2) \rightarrow \text{paritynat}(1_2)$.

When expressed explicitly as a mutual inductive definition, the induction principle for even and odd says that given $C : \text{even} \rightarrow \mathcal{U}$ and $D : \text{odd} \rightarrow \mathcal{U}$, along with

- $c_0 : C(0)$,
- $c_s : \prod_{(n:\text{odd})} D(n) \rightarrow C(\text{esucc}(n))$,
- $d_s : \prod_{(n:\text{even})} C(n) \rightarrow D(\text{osucc}(n))$,

there exist $f : \prod_{(n:\text{even})} C(n)$ and $g : \prod_{(n:\text{odd})} D(n)$ such that

$$\begin{aligned} f(0) &\equiv c_0 \\ f(\text{esucc}(n)) &\equiv c_s(g(n)) \\ g(\text{osucc}(n)) &\equiv d_s(f(n)). \end{aligned}$$

In particular, just as we can only induct over an inductive family “all at once”, we have to induct on even and odd simultaneously. We will not have much use for mutual inductive definitions in this book either.

A further, more radical, generalization is to allow definition of a type family $B : A \rightarrow \mathcal{U}$ in which not only the types $B(a)$, but the type A itself, is defined as part of one big induction. In other words, not only do we specify constructors for the $B(a)$ s which can take inputs from other

$B(a')$ s, as with inductive families, we also at the same time specify constructors for A itself, which can take inputs from the $B(a)$ s. This can be regarded as an inductive family in which the indices are inductively defined simultaneously with the indexed types, or as a mutual inductive definition in which one of the types can depend on the other. More complicated dependency structures are also possible. In general, these are called **inductive-inductive definitions**. For the most part, we will not use them in this book, but their higher variant (see Chapter 6) will appear in a couple of experimental examples in ??.

The last generalization we wish to mention is **inductive-recursive definitions**, in which a type is defined inductively at the same time as a *recursive* function on it. That is, we fix a known type P , and give constructors for an inductive type A and at the same time define a function $f : A \rightarrow P$ using the recursion principle for A resulting from its constructors — with the twist that the constructors of A are allowed to refer also to the values of f . We do not yet know how to justify such definitions from a homotopical perspective, and we will not use any of them in this book.

5.8 Identity types and identity systems

We now wish to point out that the *identity types*, which play so central a role in homotopy type theory, may also be considered to be defined inductively. Specifically, they are an “inductive family” with indices, in the sense of §5.7. In fact, there are *two* ways to describe identity types as an inductive family, resulting in the two induction principles described in Chapter 1, path induction and based path induction.

In both definitions, the type A is a parameter. For the first definition, we inductively define a family $=_A : A \rightarrow A \rightarrow \mathcal{U}$, with two indices belonging to A , by the following constructor:

- for any $a : A$, an element $\text{refl}_A : a =_A a$.

By analogy with the other inductive families, we may extract the induction principle from this definition. It states that given any $C : \prod_{(a,b:A)} (a =_A b) \rightarrow \mathcal{U}$, along with $d : \prod_{(a:A)} C(a, a, \text{refl}_a)$, there exists $f : \prod_{(a,b:A)} \prod_{(p:a=_A b)} C(a, b, p)$ such that $f(a, a, \text{refl}_a) \equiv d(a)$. This is exactly the path induction principle for identity types.

For the second definition, we consider one element $a_0 : A$ to be a parameter along with $A : \mathcal{U}$, and we inductively define a family $(a_0 =_A -) : A \rightarrow \mathcal{U}$, with *one* index belonging to A , by the following constructor:

- an element $\text{refl}_{a_0} : a_0 =_A a_0$.

Note that because $a_0 : A$ was fixed as a parameter, the constructor refl_{a_0} does not appear inside the inductive definition as a function, but only as an element. The induction principle for this definition says that given $C : \prod_{(b:A)} (a_0 =_A b) \rightarrow \mathcal{U}$ along with an element $d : C(a_0, \text{refl}_{a_0})$, there exists $f : \prod_{(b:A)} \prod_{(p:a_0=_A b)} C(b, p)$ with $f(a_0, \text{refl}_{a_0}) \equiv d$. This is exactly the based path induction principle for identity types.

The view of identity types as inductive types has historically caused some confusion, because of the intuition mentioned in §5.1 that all the elements of an inductive type should be obtained by repeatedly applying its constructors. For ordinary inductive types such as $\mathbf{2}$ and \mathbb{N} , this is the case: we saw in Eq. (1.8.1) that indeed every element of $\mathbf{2}$ is either 0_2 or 1_2 , and similarly one can prove that every element of \mathbb{N} is either 0 or a successor.

However, this is *not* true for identity types: there is only one constructor refl , but not every path is equal to the constant path. More precisely, we cannot prove, using only the induction

principle for identity types (either one), that every inhabitant of $a =_A a$ is equal to refl_a . In order to actually exhibit a counterexample, we need some additional principle such as the univalence axiom — recall that in Example 3.1.9 we used univalence to exhibit a particular path $\mathbf{2} =_{\mathcal{U}} \mathbf{2}$ which is not equal to refl_2 .

The point is that, as validated by the study of homotopy-initial algebras, an inductive definition should be regarded as *freely generated* by its constructors. Of course, a freely generated structure may contain elements other than its generators: for instance, the free group on two symbols x and y contains not only x and y but also words such as xy , $yx^{-1}y$, and $x^3y^2x^{-2}yx$. In general, the elements of a free structure are obtained by applying not only the generators, but also the operations of the ambient structure, such as the group operations if we are talking about free groups.

In the case of inductive types, we are talking about freely generated *types* — so what are the “operations” of the structure of a type? If types are viewed as like *sets*, as was traditionally the case in type theory, then there are no such operations, and hence we expect there to be no elements in an inductive type other than those resulting from its constructors. In homotopy type theory, we view types as like *spaces* or ∞ -groupoids, in which case there are many operations on the *paths* (concatenation, inversion, etc.) — this will be important in Chapter 6 — but there are still no operations on the *objects* (elements). Thus, it is still true for us that, e.g., every element of $\mathbf{2}$ is either 0_2 or 1_2 , and every element of \mathbb{N} is either 0 or a successor.

However, as we saw in Chapter 2, viewing types as ∞ -groupoids entails also viewing functions as functors, and this includes type families $B : A \rightarrow \mathcal{U}$. Thus, the identity type $(a_0 =_A -)$, viewed as an inductive type family, is actually a *freely generated functor* $A \rightarrow \mathcal{U}$. Specifically, it is the functor $F : A \rightarrow \mathcal{U}$ freely generated by one element $\text{refl}_{a_0} : F(a_0)$. And a functor does have operations on objects, namely the action of the morphisms (paths) of A .

In category theory, the *Yoneda lemma* tells us that for any category A and object a_0 , the functor freely generated by an element of $F(a_0)$ is the representable functor $\text{hom}_A(a_0, -)$. Thus, we should expect the identity type $(a_0 =_A -)$ to be this representable functor, and this is indeed exactly how we view it: $(a_0 =_A b)$ is the space of morphisms (paths) in A from a_0 to b .

One reason for viewing identity types as inductive families is to apply the uniqueness principles of §§5.2 and 5.5. Specifically, we can characterize the family of identity types of a type A , up to equivalence, by giving another family of types over $A \times A$ satisfying the same induction principle. This suggests the following definitions and theorem.

Definition 5.8.1. Let A be a type and $a_0 : A$ an element.

- A **pointed predicate** over (A, a_0) is a family $R : A \rightarrow \mathcal{U}$ equipped with an element $r_0 : R(a_0)$.
- For pointed predicates (R, r_0) and (S, s_0) , a family of maps $g : \prod_{(b:A)} R(b) \rightarrow S(b)$ is **pointed** if $g(a_0, r_0) = s_0$. We have

$$\text{ppmap}(R, S) := \sum_{g: \prod_{(b:A)} R(b) \rightarrow S(b)} (g(a_0, r_0) = s_0).$$

- An **identity system at a_0** is a pointed predicate (R, r_0) such that for any type family $D : \prod_{(b:A)} R(b) \rightarrow \mathcal{U}$ and $d : D(a_0, r_0)$, there exists a function $f : \prod_{(b:A)} \prod_{(r:R(b))} D(b, r)$ such that $f(a_0, r_0) = d$.

Theorem 5.8.2. For a pointed predicate (R, r_0) over (A, a_0) , the following are logically equivalent.

- (i) (R, r_0) is an identity system at a_0 .
- (ii) For any pointed predicate (S, s_0) , the type $\text{ppmap}(R, S)$ is contractible.
- (iii) For any $b : A$, the function $\text{transport}^R(-, r_0) : (a_0 =_A b) \rightarrow R(b)$ is an equivalence.
- (iv) The type $\sum_{(b:A)} R(b)$ is contractible.

Note that the equivalences (i) \Leftrightarrow (ii) \Leftrightarrow (iii) are a version of Lemma 5.5.4 for identity types $a_0 =_A -$, regarded as inductive families varying over one element of A . Of course, (ii)–(iv) are mere propositions, so that logical equivalence implies actual equivalence. (Condition (i) is also a mere proposition, but we will not prove this.) Note also that unlike (i)–(iii), statement (iv) doesn't refer to a_0 or r_0 .

Proof. First, assume (i) and let (S, s_0) be a pointed predicate. Define $D(b, r) := S(b)$ and $d := s_0 : S(a_0) \equiv D(a_0, r_0)$. Since R is an identity system, we have $f : \prod_{(b:A)} R(b) \rightarrow S(b)$ with $f(a_0, r_0) = s_0$; hence $\text{ppmap}(R, S)$ is inhabited. Now suppose $(f, f_r), (g, g_r) : \text{ppmap}(R, S)$, and define $D(b, r) := (f(b, r) = g(b, r))$, and let $d := f_r \bullet g_r^{-1} : f(a_0, r_0) = s_0 = g(a_0, r_0)$. Then again since R is an identity system, we have $h : \prod_{(b:A)} \prod_{(r:R(b))} D(b, r)$ such that $h(a_0, r_0) = f_r \bullet g_r^{-1}$. By the characterization of paths in Σ -types and path types, these data yield an equality $(f, f_r) = (g, g_r)$. Hence $\text{ppmap}(R, S)$ is an inhabited mere proposition, and thus contractible; so (ii) holds.

Now suppose (ii), and define $S(b) := (a_0 = b)$ with $s_0 := \text{refl}_{a_0} : S(a_0)$. Then (S, s_0) is a pointed predicate, and $\lambda b. \lambda p. \text{transport}^R(p, r) : \prod_{(b:A)} S(b) \rightarrow R(b)$ is a pointed family of maps from S to R . By assumption, $\text{ppmap}(R, S)$ is contractible, hence inhabited, so there also exists a pointed family of maps from R to S . And the composites in either direction are pointed families of maps from R to R and from S to S , respectively, hence equal to identities since $\text{ppmap}(R, R)$ and $\text{ppmap}(S, S)$ are contractible. Thus (iii) holds.

Now supposing (iii), condition (iv) follows from Lemma 3.11.8, using the fact that Σ -types respect equivalences (the “if” direction of Theorem 4.7.7).

Finally, assume (iv), and let $D : \prod_{(b:A)} R(b) \rightarrow \mathcal{U}$ and $d : D(a_0, r_0)$. We can equivalently express D as a family $D' : (\sum_{(b:A)} R(b)) \rightarrow \mathcal{U}$. Now since $\sum_{(b:A)} R(b)$ is contractible, we have

$$p : \prod_{u:\sum_{(b:A)} R(b)} (a_0, r_0) = u.$$

Moreover, since the path types of a contractible type are again contractible, we have $p((a_0, r_0)) = \text{refl}_{(a_0, r_0)}$. Define $f(u) := \text{transport}^{D'}(p(u), d)$, yielding $f : \prod_{(u:\sum_{(b:A)} R(b))} D'(u)$, or equivalently $f : \prod_{(b:A)} \prod_{(r:R(b))} D(b, r)$. Finally, we have

$$f(a_0, r_0) \equiv \text{transport}^{D'}(p((a_0, r_0)), d) = \text{transport}^{D'}(\text{refl}_{(a_0, r_0)}, d) = d.$$

Thus, (i) holds. \square

We can deduce a similar result for identity types $=_A$, regarded as a family varying over two elements of A .

Definition 5.8.3. An **identity system** over a type A is a family $R : A \rightarrow A \rightarrow \mathcal{U}$ equipped with a function $r_0 : \prod_{(a:A)} R(a, a)$ such that for any type family $D : \prod_{(a,b:A)} R(a, b) \rightarrow \mathcal{U}$ and $d : \prod_{(a:A)} D(a, a, r_0(a))$, there exists a function $f : \prod_{(a,b:A)} \prod_{(r:R(a,b))} D(a, b, r)$ such that $f(a, a, r_0(a)) = d(a)$ for all $a : A$.

Theorem 5.8.4. For $R : A \rightarrow A \rightarrow \mathcal{U}$ equipped with $r_0 : \prod_{(a:A)} R(a, a)$, the following are logically equivalent.

- (i) (R, r_0) is an identity system over A .
- (ii) For all $a_0 : A$, the pointed predicate $(R(a_0), r_0(a_0))$ is an identity system at a_0 .
- (iii) For any $S : A \rightarrow A \rightarrow \mathcal{U}$ and $s_0 : \prod_{(a:A)} S(a, a)$, the type

$$\sum_{(g:\prod_{(a,b:A)} R(a,b) \rightarrow S(a,b))} \prod_{(a:A)} g(a, a, r_0(a)) = s_0(a)$$

is contractible.

- (iv) For any $a, b : A$, the map $\text{transport}^{R(a)}(-, r_0(a)) : (a =_A b) \rightarrow R(a, b)$ is an equivalence.
- (v) For any $a : A$, the type $\sum_{(b:A)} R(a, b)$ is contractible.

Proof. The equivalence (i) \Leftrightarrow (ii) follows exactly the proof of equivalence between the path induction and based path induction principles for identity types; see §1.12. The equivalence with (iv) and (v) then follows from Theorem 5.8.2, while (iii) is straightforward. \square

One reason this characterization is interesting is that it provides an alternative way to state univalence and function extensionality. The univalence axiom for a universe \mathcal{U} says exactly that the type family

$$(- \simeq -) : \mathcal{U} \rightarrow \mathcal{U} \rightarrow \mathcal{U}$$

together with $\text{id} : \prod_{(A:\mathcal{U})} (A \simeq A)$ satisfies Theorem 5.8.4(iv). Therefore, it is equivalent to the corresponding version of (i), which we can state as follows.

Corollary 5.8.5 (Equivalence induction). *Given any type family $D : \prod_{(A,B:\mathcal{U})} (A \simeq B) \rightarrow \mathcal{U}$ and function $d : \prod_{(A:\mathcal{U})} D(A, A, \text{id}_A)$, there exists $f : \prod_{(A,B:\mathcal{U})} \prod_{(e:A \simeq B)} D(A, B, e)$ such that $f(A, A, \text{id}_A) = d(A)$ for all $A : \mathcal{U}$.*

In other words, to prove something about all equivalences, it suffices to prove it about identity maps. We have already used this principle (without stating it in generality) in Lemma 4.1.1.

Similarly, function extensionality says that for any $B : A \rightarrow \mathcal{U}$, the type family

$$(- \sim -) : \left(\prod_{a:A} B(a) \right) \rightarrow \left(\prod_{a:A} B(a) \right) \rightarrow \mathcal{U}$$

together with $\lambda f. \lambda a. \text{refl}_{f(a)}$ satisfies Theorem 5.8.4(iv). Thus, it is also equivalent to the corresponding version of (i).

Corollary 5.8.6 (Homotopy induction). *Given any $D : \prod_{(f,g:\prod_{(a:A)} B(a))} (f \sim g) \rightarrow \mathcal{U}$ and $d : \prod_{(f:\prod_{(a:A)} B(a))} D(f, f, \lambda x. \text{refl}_{f(x)})$, there exists*

$$k : \prod_{(f,g:\prod_{(a:A)} B(a))} \prod_{(h:f \sim g)} D(f, g, h)$$

such that $k(f, f, \lambda x. \text{refl}_{f(x)}) = d(f)$ for all f .

Notes

Inductive definitions have a long pedigree in mathematics, arguably going back at least to Frege and Peano's axioms for the natural numbers. More general “inductive predicates” are not uncommon, but in set theoretic foundations they are usually constructed explicitly, either as an

intersection of an appropriate class of subsets or using transfinite iteration along the ordinals, rather than regarded as a basic notion.

In type theory, particular cases of inductive definitions date back to Martin-Löf's original papers: [ML71] presents a general notion of inductively defined predicates and relations; the notion of inductive type was present (but only with instances, not as a general notion) in Martin-Löf's first papers in type theory [ML75]; and then as a general notion with W-types in [ML82].

A general notion of inductive type was introduced in 1985 by Constable and Mendler [CM85]. A general schema for inductive types in intensional type theory was suggested in [PPM90]. Further developments included [CP90, Dyb91].

The notion of inductive-recursive definition appears in [Dyb00]. An important type-theoretic notion is the notion of tree types (a general expression of the notion of Post system in type theory) which appears in [PS89].

The universal property of the natural numbers as an initial object of the category of \mathbb{N} -algebras is due to Lawvere [Law06]. This was later generalized to a description of W-types as initial algebras for polynomial endofunctors by [MP00]. The coherently homotopy-theoretic equivalence between such universal properties and the corresponding induction principles (§§5.4 and 5.5) is due to [AGS12].

For actual constructions of inductive types in homotopy-theoretic semantics of type theory, see [KLV12, vdBM15, LS17].

Exercises

Exercise 5.1. Derive the induction principle for the type $\text{List}(A)$ of lists from its definition as an inductive type in §5.1.

Exercise 5.2. Construct two functions on natural numbers which satisfy the same recurrence (e_z, e_s) judgmentally, but are not judgmentally equal.

Exercise 5.3. Construct two different recurrences (e_z, e_s) on the same type E which are both satisfied judgmentally by the same function $f : \mathbb{N} \rightarrow E$.

Exercise 5.4. Show that for any type family $E : \mathbf{2} \rightarrow \mathcal{U}$, the induction operator

$$\text{ind}_2(E) : (E(0_2) \times E(1_2)) \rightarrow \prod_{b:\mathbf{2}} E(b)$$

is an equivalence.

Exercise 5.5. Show that the analogous statement to Exercise 5.4 for \mathbb{N} fails.

Exercise 5.6. Show that if we assume simple instead of dependent elimination for W-types, the uniqueness property (analogue of Theorem 5.3.1) fails to hold. That is, exhibit a type satisfying the recursion principle of a W-type, but for which functions are not determined uniquely by their recurrence.

Exercise 5.7. Suppose that in the “inductive definition” of the type C at the beginning of §5.6, we replace the type \mathbb{N} by $\mathbf{0}$. Analogously to (5.6.1), we might consider a recursion principle for this type with hypothesis

$$h : (C \rightarrow \mathbf{0}) \rightarrow (P \rightarrow \mathbf{0}) \rightarrow P.$$

Show that even without a computation rule, this recursion principle is inconsistent, i.e. it allows us to construct an element of $\mathbf{0}$.

Exercise 5.8. Consider now an “inductive type” D with one constructor $\text{scott} : (D \rightarrow D) \rightarrow D$. The second recursor for C suggested in §5.6 leads to the following recursor for D :

$$\text{rec}_D : \prod_{P:\mathcal{U}} ((D \rightarrow D) \rightarrow (D \rightarrow P) \rightarrow P) \rightarrow D \rightarrow P$$

with computation rule $\text{rec}_D(P, h, \text{scott}(\alpha)) \equiv h(\alpha, (\lambda d. \text{rec}_D(P, h, \alpha(d))))$. Show that this also leads to a contradiction.

Exercise 5.9. Let A be an arbitrary type and consider generally an “inductive definition” of a type L_A with constructor $\text{lawvere} : (L_A \rightarrow A) \rightarrow L_A$. The second recursor for C suggested in §5.6 leads to the following recursor for L_A :

$$\text{rec}_{L_A} : \prod_{P:\mathcal{U}} ((L_A \rightarrow A) \rightarrow P) \rightarrow L_A \rightarrow P$$

with computation rule $\text{rec}_{L_A}(P, h, \text{lawvere}(\alpha)) \equiv h(\alpha)$. Using this, show that A has the **fixed-point property**, i.e. for every function $f : A \rightarrow A$ there exists an $a : A$ such that $f(a) = a$. In particular, L_A is inconsistent if A is a type without the fixed-point property, such as $\mathbf{0}$, $\mathbf{2}$, or \mathbb{N} .

Exercise 5.10. Continuing from Exercise 5.9, consider L_1 , which is not obviously inconsistent since $\mathbf{1}$ does have the fixed-point property. Formulate an induction principle for L_1 and its computation rule, analogously to its recursor, and using this, prove that it is contractible.

Exercise 5.11. In §5.1 we defined the type $\text{List}(A)$ of finite lists of elements of some type A . Consider a similar inductive definition of a type $\text{Lost}(A)$ whose only constructor is

$$\text{cons} : A \rightarrow \text{Lost}(A) \rightarrow \text{Lost}(A).$$

Show that $\text{Lost}(A)$ is equivalent to $\mathbf{0}$.

Exercise 5.12. Suppose A is a mere proposition, and $B : A \rightarrow \mathcal{U}$.

- (i) Show that $\mathsf{W}_{(a:A)} B(a)$ is a mere proposition.
- (ii) Show that $\mathsf{W}_{(a:A)} B(a)$ is equivalent to $\sum_{(a:A)} \neg B(a)$.
- (iii) Without using $\mathsf{W}_{(a:A)} B(a)$, show that $\sum_{(a:A)} \neg B(a)$ is a homotopy W -type $\mathsf{W}_{(a:A)}^h B(a)$ in the sense of §5.5.

Exercise 5.13. Let $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$.

- (i) Show that $(\sum_{(a:A)} \neg B(a)) \rightarrow (\mathsf{W}_{(a:A)} B(a))$.
- (ii) Show that $(\mathsf{W}_{(a:A)} B(a)) \rightarrow (\neg \prod_{(a:A)} B(a))$.

Exercise 5.14. Let $A : \mathcal{U}$ and suppose that $B : A \rightarrow \mathcal{U}$ is decidable, i.e. $\prod_{(a:A)} (B(a) + \neg B(a))$ (see Definition 3.4.3). Show that $(\mathsf{W}_{(a:A)} B(a)) \rightarrow (\sum_{(a:A)} \neg B(a))$.

Exercise 5.15. Show that the following are logically equivalent.

- (i) $(\mathsf{W}_{(a:A)} B(a)) \rightarrow \|\sum_{(a:A)} \neg B(a)\|$ for any $A : \text{Set}$ and $B : A \rightarrow \text{Prop}$.
- (ii) $(\neg \prod_{(a:A)} B(a)) \rightarrow \|\mathsf{W}_{(a:A)} B(a)\|$ for any $A : \text{Set}$ and $B : A \rightarrow \text{Prop}$.
- (iii) The law of excluded middle (as in §3.4).

Similarly, using Corollary 3.2.7, show that it is inconsistent to assume that either implication in (i) or (ii) holds for all $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$.

Exercise 5.16. For $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, define

$$W'_{A,B} := \prod_{R:\mathcal{U}} \left(\prod_{a:A} (B(a) \rightarrow R) \rightarrow R \right)$$

$W'_{A,B}$ is called the **impredicative encoding** of $W_{(a:A)} B(a)$. Note that unlike $W_{(a:A)} B(a)$, it lives in a higher universe than A and B .

- (i) Show that $W'_{A,B}$ is logically equivalent (as defined in §1.11) to $W_{(a:A)} B(a)$.
- (ii) Show that $W'_{A,B}$ implies $\neg\neg \sum_{(a:A)} \neg B(a)$.
- (iii) Without using $W_{(a:A)} B(a)$, show that $W'_{A,B}$ satisfies the same *recursion principle* as $W_{(a:A)} B(a)$ for defining functions into types in the universe \mathcal{U} (to which it itself does not belong).
- (iv) Using LEM, give an example of an $A : \mathcal{U}$ and a $B : A \rightarrow \mathcal{U}$ such that $W'_{A,B}$ is not equivalent to $W_{(a:A)} B(a)$.

Exercise 5.17. Show that for any $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$, we have

$$\neg(W_{(a:A)} B(a)) \simeq \neg\left(\sum_{a:A} \neg B(a)\right).$$

In other words, $W_{(a:A)} B(a)$ is empty if and only if it has no nullary constructor. (Compare to Exercise 5.11.)

Chapter 6

Higher inductive types

6.1 Introduction

Like the general inductive types we discussed in Chapter 5, *higher inductive types* are a general schema for defining new types generated by some constructors. But unlike ordinary inductive types, in defining a higher inductive type we may have “constructors” which generate not only *points* of that type, but also *paths* and higher paths in that type. For instance, we can consider the higher inductive type S^1 generated by

- A point base : S^1 , and
- A path loop : $\text{base} =_{S^1} \text{base}$.

This should be regarded as entirely analogous to the definition of, for instance, $\mathbf{2}$, as being generated by

- A point $0_2 : \mathbf{2}$ and
- A point $1_2 : \mathbf{2}$,

or the definition of \mathbb{N} as generated by

- A point $0 : \mathbb{N}$ and
- A function $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$.

When we think of types as higher groupoids, the more general notion of “generation” is very natural: since a higher groupoid is a “multi-sorted object” with paths and higher paths as well as points, we should allow “generators” in all dimensions.

We will refer to the ordinary sort of constructors (such as `base`) as **point constructors** or *ordinary constructors*, and to the others (such as `loop`) as **path constructors** or *higher constructors*. Each path constructor must specify the starting and ending point of the path, which we call its **source** and **target**; for `loop`, both source and target are `base`.

Note that a path constructor such as `loop` generates a *new* inhabitant of an identity type, which is not (at least, not *a priori*) equal to any previously existing such inhabitant. In particular, `loop` is not *a priori* equal to `reflbase` (although proving that they are definitely unequal takes a little thought; see Lemma 6.4.1). This is what distinguishes S^1 from the ordinary inductive type $\mathbf{1}$.

There are some important points to be made regarding this generalization.

First of all, the word “generation” should be taken seriously, in the same sense that a group can be freely generated by some set. In particular, because a higher groupoid comes with *operations* on paths and higher paths, when such an object is “generated” by certain constructors, the operations create more paths that do not come directly from the constructors themselves. For instance, in the higher inductive type S^1 , the constructor loop is not the only nontrivial path from base to base; we have also “ $\text{loop} \cdot \text{loop}$ ” and “ $\text{loop} \cdot \text{loop} \cdot \text{loop}$ ” and so on, as well as loop^{-1} , etc., all of which are different. This may seem so obvious as to be not worth mentioning, but it is a departure from the behavior of “ordinary” inductive types, where one can expect to see nothing in the inductive type except what was “put in” directly by the constructors.

Secondly, this generation is really *free* generation: higher inductive types do not technically allow us to impose “axioms”, such as forcing “ $\text{loop} \cdot \text{loop}$ ” to equal $\text{refl}_{\text{base}}$. However, in the world of ∞ -groupoids, there is little difference between “free generation” and “presentation”, since we can make two paths equal *up to homotopy* by adding a new 2-dimensional generator relating them (e.g. a path $\text{loop} \cdot \text{loop} = \text{refl}_{\text{base}}$ in $\text{base} = \text{base}$). We do then, of course, have to worry about whether this new generator should satisfy its own “axioms”, and so on, but in principle any “presentation” can be transformed into a “free” one by making axioms into constructors. As we will see, by adding “truncation constructors” we can use higher inductive types to express classical notions such as group presentations as well.

Thirdly, even though a higher inductive type contains “constructors” which generate *paths in* that type, it is still an inductive definition of a *single* type. In particular, as we will see, it is the higher inductive type itself which is given a universal property (expressed, as usual, by an induction principle), and *not* its identity types. The identity type of a higher inductive type retains the usual induction principle of any identity type (i.e. path induction), and does not acquire any new induction principle.

Thus, it may be nontrivial to identify the identity types of a higher inductive type in a concrete way, in contrast to how in Chapter 2 we were able to give explicit descriptions of the behavior of identity types under all the traditional type forming operations. For instance, are there any paths from base to base in S^1 which are not simply composites of copies of loop and its inverse? Intuitively, it seems that the answer should be no (and it is), but proving this is not trivial. Indeed, such questions bring us rapidly to problems such as calculating the homotopy groups of spheres, a long-standing problem in algebraic topology for which no simple formula is known. Homotopy type theory brings a new and powerful viewpoint to bear on such questions, but it also requires type theory to become as complex as the answers to these questions.

Fourthly, the “dimension” of the constructors (i.e. whether they output points, paths, paths between paths, etc.) does not have a direct connection to which dimensions the resulting type has nontrivial homotopy in. As a simple example, if an inductive type B has a constructor of type $A \rightarrow B$, then any paths and higher paths in A result in paths and higher paths in B , even though the constructor is not a “higher” constructor at all. The same thing happens with higher constructors too: having a constructor of type $A \rightarrow (x =_B y)$ means not only that points of A yield paths from x to y in B , but that paths in A yield paths between these paths, and so on. As we will see, this possibility is responsible for much of the power of higher inductive types.

On the other hand, it is even possible for constructors *without* higher types in their inputs to generate “unexpected” higher paths. For instance, in the 2-dimensional sphere S^2 generated by

- A point $\text{base} : S^2$, and
- A 2-dimensional path $\text{surf} : \text{refl}_{\text{base}} = \text{refl}_{\text{base}}$ in $\text{base} = \text{base}$,

there is a nontrivial 3-dimensional path from $\text{refl}_{\text{refl}_{\text{base}}}$ to itself. Topologists will recognize this

path as an incarnation of the *Hopf fibration*. From a category-theoretic point of view, this is the same sort of phenomenon as the fact mentioned above that S^1 contains not only loop but also loop • loop and so on: it's just that in a *higher* groupoid, there are *operations* which raise dimension. Indeed, we saw many of these operations back in §2.1: the associativity and unit laws are not just properties, but operations, whose inputs are 1-paths and whose outputs are 2-paths.

6.2 Induction principles and dependent paths

When we describe a higher inductive type such as the circle as being generated by certain constructors, we have to explain what this means by giving rules analogous to those for the basic type constructors from Chapter 1. The constructors themselves give the *introduction* rules, but it requires a bit more thought to explain the *elimination* rules, i.e. the induction and recursion principles. In this book we do not attempt to give a general formulation of what constitutes a “higher inductive definition” and how to extract the elimination rule from such a definition — indeed, this is a subtle question and the subject of current research. Instead we will rely on some general informal discussion and numerous examples.

The recursion principle is usually easy to describe: given any type equipped with the same structure with which the constructors equip the higher inductive type in question, there is a function which maps the constructors to that structure. For instance, in the case of S^1 , the recursion principle says that given any type B equipped with a point $b : B$ and a path $\ell : b = b$, there is a function $f : S^1 \rightarrow B$ such that $f(\text{base}) = b$ and $\text{ap}_f(\text{loop}) = \ell$.

The latter two equalities are the *computation rules*. There is, however, a question of whether these computation rules are judgmental equalities or propositional equalities (paths). For ordinary inductive types, we had no qualms about making them judgmental, although we saw in Chapter 5 that making them propositional would still yield the same type up to equivalence. In the ordinary case, one may argue that the computation rules are really *definitional* equalities, in the intuitive sense described in the Introduction.

For higher inductive types, this is less clear. Moreover, since the operation ap_f is not really a fundamental part of the type theory, but something that we *defined* using the induction principle of identity types (and which we might have defined in some other, equivalent, way), it seems inappropriate to refer to it explicitly in a *judgmental* equality. Judgmental equalities are part of the deductive system, which should not depend on particular choices of definitions that we may make *within* that system. There are also semantic and implementation issues to consider; see the Notes.

It does seem unproblematic to make the computational rules for the *point* constructors of a higher inductive type judgmental. In the example above, this means we have $f(\text{base}) \equiv b$, judgmentally. This choice facilitates a computational view of higher inductive types. Moreover, it also greatly simplifies our lives, since otherwise the second computation rule $\text{ap}_f(\text{loop}) = \ell$ would not even be well-typed as a propositional equality; we would have to compose one side or the other with the specified identification of $f(\text{base})$ with b . (Such problems do arise eventually, of course, when we come to talk about paths of higher dimension, but that will not be of great

concern to us here. See also §6.7.) Thus, we take the computation rules for point constructors to be judgmental, and those for paths and higher paths to be propositional.¹

Remark 6.2.1. Recall that for ordinary inductive types, we regard the computation rules for a recursively defined function as not merely judgmental equalities, but *definitional* ones, and thus we may use the notation \equiv for them. For instance, the truncated predecessor function $p : \mathbb{N} \rightarrow \mathbb{N}$ is defined by $p(0) \equiv 0$ and $p(\text{succ}(n)) \equiv n$. In the case of higher inductive types, this sort of notation is reasonable for the point constructors (e.g. $f(\text{base}) \equiv b$), but for the path constructors it could be misleading, since equalities such as $f(\text{loop}) = \ell$ are not judgmental. Thus, we hybridize the notations, writing instead $f(\text{loop}) := \ell$ for this sort of “propositional equality by definition”.

Now, what about the induction principle (the dependent eliminator)? Recall that for an ordinary inductive type W , to prove by induction that $\prod_{(x:W)} P(x)$, we must specify, for each constructor of W , an operation on P which acts on the “fibers” above that constructor in W . For instance, if W is the natural numbers \mathbb{N} , then to prove by induction that $\prod_{(x:\mathbb{N})} P(x)$, we must specify

- An element $b : P(0)$ in the fiber over the constructor $0 : \mathbb{N}$, and
- For each $n : \mathbb{N}$, a function $P(n) \rightarrow P(\text{succ}(n))$.

The second can be viewed as a function “ $P \rightarrow P$ ” lying *over* the constructor $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$, generalizing how $b : P(0)$ lies over the constructor $0 : \mathbb{N}$.

By analogy, therefore, to prove that $\prod_{(x:S^1)} P(x)$, we should specify

- An element $b : P(\text{base})$ in the fiber over the constructor $\text{base} : S^1$, and
- A path from b to b “lying over the constructor $\text{loop} : \text{base} = \text{base}$ ”.

Note that even though S^1 contains paths other than loop (such as $\text{refl}_{\text{base}}$ and $\text{loop} \cdot \text{loop}$), we only need to specify a path lying over the constructor *itself*. This expresses the intuition that S^1 is “freely generated” by its constructors.

The question, however, is what it means to have a path “lying over” another path. It definitely does *not* mean simply a path $b = b$, since that would be a path in the fiber $P(\text{base})$ (topologically, a path lying over the *constant* path at base). Actually, however, we have already answered this question in Chapter 2: in the discussion preceding Lemma 2.3.4 we concluded that a path from $u : P(x)$ to $v : P(y)$ lying over $p : x = y$ can be represented by a path $p_*(u) = v$ in the fiber $P(y)$. Since we will have a lot of use for such **dependent paths** in this chapter, we introduce a special notation for them:

$$(u =_p^P v) : \equiv (\text{transport}^P(p, u) = v). \quad (6.2.2)$$

Remark 6.2.3. There are other possible ways to define dependent paths. For instance, instead of $p_*(u) = v$ we could consider $u = (p^{-1})_*(v)$. We could also obtain it as a special case of a more general “heterogeneous equality”, or with a direct definition as an inductive type family. All these definitions result in equivalent types, so in that sense it doesn’t much matter which we

¹In particular, in the language of §1.1, this means that our higher inductive types are a mix of *rules* (specifying how we can introduce such types and their elements, their induction principle, and their computation rules for point constructors) and *axioms* (the computation rules for path constructors, which assert that certain identity types are inhabited by otherwise unspecified terms). We may hope that eventually, there will be a better type theory in which higher inductive types, like univalence, will be presented using only rules and no axioms.

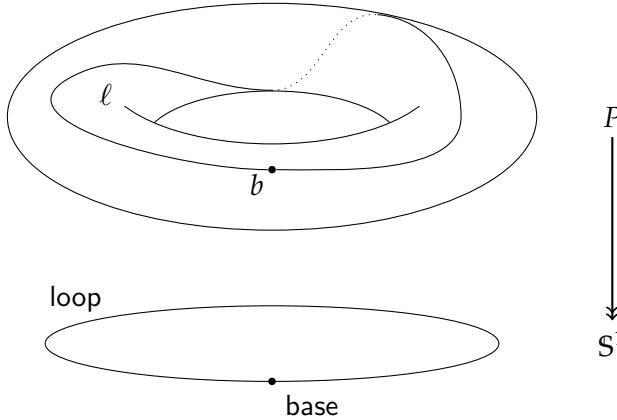


Figure 6.1: The topological induction principle for S^1

pick. However, choosing $p_*(u) = v$ as the definition makes it easiest to conclude other things about dependent paths, such as the fact that apd_f produces them, or that we can compute them in particular type families using the transport lemmas in §2.5.

With the notion of dependent paths in hand, we can now state more precisely the induction principle for S^1 : given $P : S^1 \rightarrow \mathcal{U}$ and

- an element $b : P(\text{base})$, and
- a path $\ell : b =_{\text{loop}}^P b$,

there is a function $f : \prod_{(x:S^1)} P(x)$ such that $f(\text{base}) \equiv b$ and $\text{apd}_f(\text{loop}) = \ell$. As in the non-dependent case, we speak of defining f by $f(\text{base}) := b$ and $\text{apd}_f(\text{loop}) := \ell$.

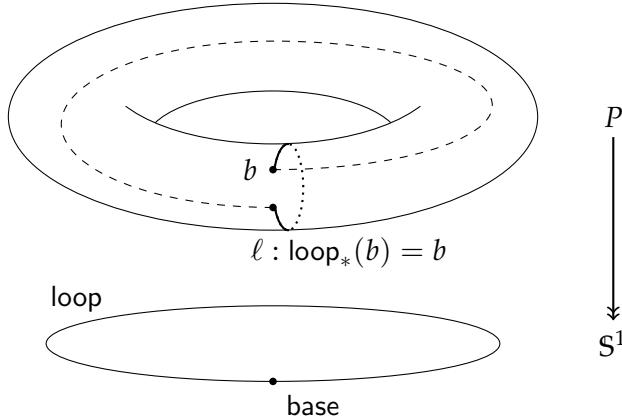
Remark 6.2.4. When describing an application of this induction principle informally, we regard it as a splitting of the goal “ $P(x)$ for all $x : S^1$ ” into two cases, which we will sometimes introduce with phrases such as “when x is base” and “when x varies along loop”, respectively. There is no specific mathematical meaning assigned to “varying along a path”: it is just a convenient way to indicate the beginning of the corresponding section of a proof; see Lemma 6.4.2 for an example.

Topologically, the induction principle for S^1 can be visualized as shown in Figure 6.1. Given a fibration over the circle (which in the picture is a torus), to define a section of this fibration is the same as to give a point b in the fiber over base along with a path from b to b lying over loop. The way we interpret this type-theoretically, using our definition of dependent paths, is shown in Figure 6.2: the path from b to b over loop is represented by a path from $\text{loop}_*(b)$ to b in the fiber over base.

Of course, we expect to be able to prove the recursion principle from the induction principle, by taking P to be a constant type family. This is in fact the case, although deriving the non-dependent computation rule for loop (which refers to ap_f) from the dependent one (which refers to apd_f) is surprisingly a little tricky.

Lemma 6.2.5. *If A is a type together with $a : A$ and $p : a =_A a$, then there is a function $f : S^1 \rightarrow A$ with*

$$\begin{aligned} f(\text{base}) &:= a \\ \text{ap}_f(\text{loop}) &:= p. \end{aligned}$$

Figure 6.2: The type-theoretic induction principle for \mathbb{S}^1

Proof. We would like to apply the induction principle of \mathbb{S}^1 to the constant type family, $(\lambda x. A) : \mathbb{S}^1 \rightarrow \mathcal{U}$. The required hypotheses for this are a point of $(\lambda x. A)(\text{base}) \equiv A$, which we have (namely $a : A$), and a dependent path in $a =_{\text{loop}}^{x \mapsto A} a$, or equivalently $\text{transport}^{x \mapsto A}(\text{loop}, a) = a$. This latter type is not the same as the type $a =_A a$ where p lives, but it is equivalent to it, because by Lemma 2.3.5 we have $\text{transportconst}_{\text{loop}}^A(a) : \text{transport}^{x \mapsto A}(\text{loop}, a) = a$. Thus, given $a : A$ and $p : a = a$, we can consider the composite

$$\text{transportconst}_{\text{loop}}^A(a) \bullet p : (a =_{\text{loop}}^{x \mapsto A} a).$$

Applying the induction principle, we obtain $f : \mathbb{S}^1 \rightarrow A$ such that

$$f(\text{base}) \equiv a \quad \text{and} \tag{6.2.6}$$

$$\text{apd}_f(\text{loop}) = \text{transportconst}_{\text{loop}}^A(a) \bullet p. \tag{6.2.7}$$

It remains to derive the equality $\text{apd}_f(\text{loop}) = p$. However, by Lemma 2.3.8, we have

$$\text{apd}_f(\text{loop}) = \text{transportconst}_{\text{loop}}^A(f(\text{base})) \bullet \text{ap}_f(\text{loop}).$$

Combining this with (6.2.7) and canceling the occurrences of transportconst (which are the same by (6.2.6)), we obtain $\text{ap}_f(\text{loop}) = p$. \square

We also have a corresponding uniqueness principle.

Lemma 6.2.8. *If A is a type and $f, g : \mathbb{S}^1 \rightarrow A$ are two maps together with two equalities p, q :*

$$\begin{aligned} p &: f(\text{base}) =_A g(\text{base}), \\ q &: f(\text{loop}) =_p^{\lambda x. x =_A x} g(\text{loop}). \end{aligned}$$

Then for all $x : \mathbb{S}^1$ we have $f(x) = g(x)$.

Proof. We apply the induction principle of \mathbb{S}^1 at the type family $P(x) := (f(x) = g(x))$. When x is base, p is exactly what we need. And when x varies along loop, we need $p =_{\text{loop}}^{\lambda x. f(x) = g(x)} p$, which by Theorems 2.11.3 and 2.11.5 can be reduced to q . \square

These two lemmas imply the expected universal property of the circle:

Lemma 6.2.9. *For any type A we have a natural equivalence*

$$(\mathbb{S}^1 \rightarrow A) \simeq \sum_{x:A} (x = x).$$

Proof. We have a canonical function $f : (\mathbb{S}^1 \rightarrow A) \rightarrow \sum_{(x:A)} (x = x)$ defined by $f(g) := (g(\text{base}), g(\text{loop}))$. The induction principle shows that the fibers of f are inhabited, while the uniqueness principle shows that they are mere propositions. Hence they are contractible, so f is an equivalence. \square

As in §5.5, we can show that the conclusion of Lemma 6.2.9 is equivalent to having an induction principle with propositional computation rules. Other higher inductive types also satisfy lemmas analogous to Lemmas 6.2.5 and 6.2.9; we will generally leave their proofs to the reader. We now proceed to consider many examples.

6.3 The interval

The **interval**, which we denote I , is perhaps an even simpler higher inductive type than the circle. It is generated by:

- a point $0_I : I$,
- a point $1_I : I$, and
- a path $\text{seg} : 0_I =_I 1_I$.

The recursion principle for the interval says that given a type B along with

- a point $b_0 : B$,
- a point $b_1 : B$, and
- a path $s : b_0 = b_1$,

there is a function $f : I \rightarrow B$ such that $f(0_I) \equiv b_0$, $f(1_I) \equiv b_1$, and $f(\text{seg}) = s$. The induction principle says that given $P : I \rightarrow \mathcal{U}$ along with

- a point $b_0 : P(0_I)$,
- a point $b_1 : P(1_I)$, and
- a path $s : b_0 =_{\text{seg}}^P b_1$,

there is a function $f : \prod_{(x:I)} P(x)$ such that $f(0_I) \equiv b_0$, $f(1_I) \equiv b_1$, and $\text{apd}_f(\text{seg}) = s$.

Regarded purely up to homotopy, the interval is not really interesting:

Lemma 6.3.1. *The type I is contractible.*

Proof. We prove that for all $x : I$ we have $x =_I 1_I$. In other words we want a function f of type $\prod_{(x:I)} (x =_I 1_I)$. We begin to define f in the following way:

$$\begin{aligned} f(0_I) &:= \text{seg} : 0_I =_I 1_I, \\ f(1_I) &:= \text{refl}_{1_I} : 1_I =_I 1_I. \end{aligned}$$

It remains to define $\text{apd}_f(\text{seg})$, which must have type $\text{seg} =_{\text{seg}}^{\lambda x. x =_I 1_I} \text{refl}_{1_I}$. By definition this type is $\text{seg}_*(\text{seg}) =_{1_I =_I 1_I} \text{refl}_{1_I}$, which in turn is equivalent to $\text{seg}^{-1} \bullet \text{seg} = \text{refl}_{1_I}$. But there is a canonical element of that type, namely the proof that path inverses are in fact inverses. \square

However, type-theoretically the interval does still have some interesting features, just like the topological interval in classical homotopy theory. For instance, it enables us to give an easy proof of function extensionality. (Of course, as in §4.9, for the duration of the following proof we suspend our overall assumption of the function extensionality axiom.)

Lemma 6.3.2. *If $f, g : A \rightarrow B$ are two functions such that $f(x) = g(x)$ for every $x : A$, then $f = g$ in the type $A \rightarrow B$.*

Proof. Let's call the proof we have $p : \prod_{(x:A)}(f(x) = g(x))$. For all $x : A$ we define a function $\tilde{p}_x : I \rightarrow B$ by

$$\begin{aligned}\tilde{p}_x(0_I) &:= f(x), \\ \tilde{p}_x(1_I) &:= g(x), \\ \tilde{p}_x(\text{seg}) &:= p(x).\end{aligned}$$

We now define $q : I \rightarrow (A \rightarrow B)$ by

$$q(i) := (\lambda x. \tilde{p}_x(i))$$

Then $q(0_I)$ is the function $\lambda x. \tilde{p}_x(0_I)$, which is equal to f because $\tilde{p}_x(0_I)$ is defined by $f(x)$. Similarly, we have $q(1_I) = g$, and hence

$$q(\text{seg}) : f =_{(A \rightarrow B)} g$$

□

In Exercise 6.10 we ask the reader to complete the proof of the full function extensionality axiom from Lemma 6.3.2.

6.4 Circles and spheres

We have already discussed the circle S^1 as the higher inductive type generated by

- A point $\text{base} : S^1$, and
- A path $\text{loop} : \text{base} =_{S^1} \text{base}$.

Its induction principle says that given $P : S^1 \rightarrow \mathcal{U}$ along with $b : P(\text{base})$ and $\ell : b =_{\text{loop}}^P b$, we have $f : \prod_{(x:S^1)} P(x)$ with $f(\text{base}) \equiv b$ and $\text{apd}_f(\text{loop}) = \ell$. Its non-dependent recursion principle says that given B with $b : B$ and $\ell : b = b$, we have $f : S^1 \rightarrow B$ with $f(\text{base}) \equiv b$ and $f(\text{loop}) = \ell$.

We observe that the circle is nontrivial.

Lemma 6.4.1. $\text{loop} \neq \text{refl}_{\text{base}}$.

Proof. Suppose that $\text{loop} = \text{refl}_{\text{base}}$. Then since for any type A with $x : A$ and $p : x = x$, there is a function $f : S^1 \rightarrow A$ defined by $f(\text{base}) := x$ and $f(\text{loop}) := p$, we have

$$p = f(\text{loop}) = f(\text{refl}_{\text{base}}) = \text{refl}_x.$$

But this implies that every type is a set, which as we have seen is not the case (see Example 3.1.9). □

The circle also has the following interesting property, which is useful as a source of counterexamples.

Lemma 6.4.2. *There exists an element of $\prod_{(x:S^1)}(x = x)$ which is not equal to $x \mapsto \text{refl}_x$.*

Proof. We define $f : \prod_{(x:S^1)}(x = x)$ by S^1 -induction. When x is base, we let $f(\text{base}) := \text{loop}$. Now when x varies along loop (see Remark 6.2.4), we must show that $\text{transport}^{x \mapsto x=x}(\text{loop}, \text{loop}) = \text{loop}$. However, in §2.11 we observed that $\text{transport}^{x \mapsto x=x}(p, q) = p^{-1} \cdot q \cdot p$, so what we have to show is that $\text{loop}^{-1} \cdot \text{loop} \cdot \text{loop} = \text{loop}$. But this is clear by canceling an inverse.

To show that $f \neq (x \mapsto \text{refl}_x)$, it suffices to show that $f(\text{base}) \neq \text{refl}_{\text{base}}$. But $f(\text{base}) = \text{loop}$, so this is just the previous lemma. \square

For instance, this enables us to extend Example 3.1.9 by showing that any universe which contains the circle cannot be a 1-type.

Corollary 6.4.3. *If the type S^1 belongs to some universe \mathcal{U} , then \mathcal{U} is not a 1-type.*

Proof. The type $S^1 = S^1$ in \mathcal{U} is, by univalence, equivalent to the type $S^1 \simeq S^1$ of autoequivalences of S^1 , so it suffices to show that $S^1 \simeq S^1$ is not a set. For this, it suffices to show that its equality type $\text{id}_{S^1} =_{(S^1 \simeq S^1)} \text{id}_{S^1}$ is not a mere proposition. Since being an equivalence is a mere proposition, this type is equivalent to $\text{id}_{S^1} =_{(S^1 \rightarrow S^1)} \text{id}_{S^1}$. But by function extensionality, this is equivalent to $\prod_{(x:S^1)}(x = x)$, which as we have seen in Lemma 6.4.2 contains two unequal elements. \square

We have also mentioned that the 2-sphere S^2 should be the higher inductive type generated by

- A point $\text{base} : S^2$, and
- A 2-dimensional path $\text{surf} : \text{refl}_{\text{base}} = \text{refl}_{\text{base}}$ in $\text{base} = \text{base}$.

The recursion principle for S^2 is not hard: it says that given B with $b : B$ and $s : \text{refl}_b = \text{refl}_b$, we have $f : S^2 \rightarrow B$ with $f(\text{base}) \equiv b$ and $\text{ap}_f^2(\text{surf}) = s$. Here by “ $\text{ap}_f^2(\text{surf})$ ” we mean an extension of the functorial action of f to two-dimensional paths, which can be stated precisely as follows.

Lemma 6.4.4. *Given $f : A \rightarrow B$ and $x, y : A$ and $p, q : x = y$, and $r : p = q$, we have a path $\text{ap}_f^2(r) : f(p) = f(q)$.*

Proof. By path induction, we may assume $p \equiv q$ and r is reflexivity. But then we may define $\text{ap}_f^2(\text{refl}_p) := \text{refl}_{f(p)}$. \square

In order to state the general induction principle, we need a version of this lemma for dependent functions, which in turn requires a notion of dependent two-dimensional paths. As before, there are many ways to define such a thing; one is by way of a two-dimensional version of transport.

Lemma 6.4.5. *Given $P : A \rightarrow \mathcal{U}$ and $x, y : A$ and $p, q : x = y$ and $r : p = q$, for any $u : P(x)$ we have $\text{transport}^2(r, u) : p_*(u) = q_*(u)$.*

Proof. By path induction. \square

Now suppose given $x, y : A$ and $p, q : x = y$ and $r : p = q$ and also points $u : P(x)$ and $v : P(y)$ and dependent paths $h : u =_p^P v$ and $k : u =_q^P v$. By our definition of dependent paths, this means $h : p_*(u) = v$ and $k : q_*(u) = v$. Thus, it is reasonable to define the type of dependent 2-paths over r to be

$$(h =_r^P k) := (h = \text{transport}^2(r, u) \cdot k).$$

We can now state the dependent version of Lemma 6.4.4.

Lemma 6.4.6. *Given $P : A \rightarrow \mathcal{U}$ and $x, y : A$ and $p, q : x = y$ and $r : p = q$ and a function $f : \prod_{(x:A)} P(x)$, we have $\text{apd}_f^2(r) : \text{apd}_f(p) =_r^P \text{apd}_f(q)$.*

Proof. Path induction. □

Now we can state the induction principle for S^2 : suppose we are given $P : S^2 \rightarrow \mathcal{U}$ with $b : P(\text{base})$ and $s : \text{refl}_b =_{\text{surf}}^Q \text{refl}_b$ where $Q := \lambda p. b =_p^P b$. Then there is a function $f : \prod_{(x:S^2)} P(x)$ such that $f(\text{base}) \equiv b$ and $\text{apd}_f^2(\text{surf}) = s$.

Of course, this explicit approach gets more and more complicated as we go up in dimension. Thus, if we want to define n -spheres for all n , we need some more systematic idea. One approach is to work with n -dimensional loops directly, rather than general n -dimensional paths.

Recall from §2.1 the definitions of *pointed types* \mathcal{U}_* , and the n -fold loop space $\Omega^n : \mathcal{U}_* \rightarrow \mathcal{U}_*$ (Definitions 2.1.7 and 2.1.8). Now we can define the n -sphere S^n to be the higher inductive type generated by

- A point $\text{base} : S^n$, and
- An n -loop $\text{loop}_n : \Omega^n(S^n, \text{base})$.

In order to write down the induction principle for this presentation, we would need to define a notion of “dependent n -loop”, along with the action of dependent functions on n -loops. We leave this to the reader (see Exercise 6.4); in the next section we will discuss a different way to define the spheres that is sometimes more tractable.

6.5 Suspensions

The **suspension** of a type A is the universal way of making the points of A into paths (and hence the paths in A into 2-paths, and so on). It is a type ΣA defined by the following generators:²

- a point $\text{N} : \Sigma A$,
- a point $\text{S} : \Sigma A$, and
- a function $\text{merid} : A \rightarrow (\text{N} =_{\Sigma A} \text{S})$.

The names are intended to suggest a “globe” of sorts, with a north pole, a south pole, and an A ’s worth of meridians from one to the other. Indeed, as we will see, if $A = S^1$, then its suspension is equivalent to the surface of an ordinary sphere, S^2 .

The recursion principle for ΣA says that given a type B together with

- points $n, s : B$ and
- a function $m : A \rightarrow (n = s)$,

we have a function $f : \Sigma A \rightarrow B$ such that $f(\text{N}) \equiv n$ and $f(\text{S}) \equiv s$, and for all $a : A$ we have $f(\text{merid}(a)) = m(a)$. Similarly, the induction principle says that given $P : \Sigma A \rightarrow \mathcal{U}$ together with

- a point $n : P(\text{N})$,
- a point $s : P(\text{S})$, and
- for each $a : A$, a path $m(a) : n =_{\text{merid}(a)}^P s$,

²There is an unfortunate clash of notation with dependent pair types, which of course are also written with a Σ . However, context usually disambiguates.

there exists a function $f : \prod_{(x:\Sigma A)} P(x)$ such that $f(\mathsf{N}) \equiv n$ and $f(\mathsf{S}) \equiv s$ and for each $a : A$ we have $\text{apd}_f(\text{merid}(a)) = m(a)$.

Our first observation about suspension is that it gives another way to define the circle.

Lemma 6.5.1. $\Sigma 2 \simeq S^1$.

Proof. Define $f : \Sigma 2 \rightarrow S^1$ by recursion such that $f(\mathsf{N}) \equiv \text{base}$ and $f(\mathsf{S}) \equiv \text{base}$, while $f(\text{merid}(0_2)) \equiv \text{loop}$ but $f(\text{merid}(1_2)) \equiv \text{refl}_{\text{base}}$. Define $g : S^1 \rightarrow \Sigma 2$ by recursion such that $g(\text{base}) \equiv \mathsf{N}$ and $g(\text{loop}) \equiv \text{merid}(0_2) \bullet \text{merid}(1_2)^{-1}$. We now show that f and g are quasi-inverses.

First we show by induction that $g(f(x)) = x$ for all $x : \Sigma 2$. If $x \equiv \mathsf{N}$, then $g(f(\mathsf{N})) \equiv g(\text{base}) \equiv \mathsf{N}$, so we have $\text{refl}_{\mathsf{N}} : g(f(\mathsf{N})) = \mathsf{N}$. If $x \equiv \mathsf{S}$, then $g(f(\mathsf{S})) \equiv g(\text{base}) \equiv \mathsf{N}$, and we choose the equality $\text{merid}(1_2) : g(f(\mathsf{S})) = \mathsf{S}$. It remains to show that for any $y : 2$, these equalities are preserved as x varies along $\text{merid}(y)$, which is to say that when refl_{N} is transported along $\text{merid}(y)$ it yields $\text{merid}(1_2)$. By transport in path spaces and pulled back fibrations, this means we are to show that

$$g(f(\text{merid}(y)))^{-1} \bullet \text{refl}_{\mathsf{N}} \bullet \text{merid}(y) = \text{merid}(1_2).$$

Of course, we may cancel refl_{N} . Now by 2-induction, we may assume either $y \equiv 0_2$ or $y \equiv 1_2$. If $y \equiv 0_2$, then we have

$$\begin{aligned} g(f(\text{merid}(0_2)))^{-1} \bullet \text{merid}(0_2) &= g(\text{loop})^{-1} \bullet \text{merid}(0_2) \\ &= (\text{merid}(0_2) \bullet \text{merid}(1_2)^{-1})^{-1} \bullet \text{merid}(0_2) \\ &= \text{merid}(1_2) \bullet \text{merid}(0_2)^{-1} \bullet \text{merid}(0_2) \\ &= \text{merid}(1_2) \end{aligned}$$

while if $y \equiv 1_2$, then we have

$$\begin{aligned} g(f(\text{merid}(1_2)))^{-1} \bullet \text{merid}(1_2) &= g(\text{refl}_{\text{base}})^{-1} \bullet \text{merid}(1_2) \\ &= \text{refl}_{\mathsf{N}}^{-1} \bullet \text{merid}(1_2) \\ &= \text{merid}(1_2). \end{aligned}$$

Thus, for all $x : \Sigma 2$, we have $g(f(x)) = x$.

Now we show by induction that $f(g(x)) = x$ for all $x : S^1$. If $x \equiv \text{base}$, then $f(g(\text{base})) \equiv f(\mathsf{N}) \equiv \text{base}$, so we have $\text{refl}_{\text{base}} : f(g(\text{base})) = \text{base}$. It remains to show that this equality is preserved as x varies along loop , which is to say that it is transported along loop to itself. Again, by transport in path spaces and pulled back fibrations, this means to show that

$$f(g(\text{loop}))^{-1} \bullet \text{refl}_{\text{base}} \bullet \text{loop} = \text{refl}_{\text{base}}.$$

However, we have

$$\begin{aligned} f(g(\text{loop})) &= f\left(\text{merid}(0_2) \bullet \text{merid}(1_2)^{-1}\right) \\ &= f(\text{merid}(0_2)) \bullet f(\text{merid}(1_2))^{-1} \\ &= \text{loop} \bullet \text{refl}_{\text{base}} \end{aligned}$$

so this follows easily. \square

Topologically, the two-point space $\mathbf{2}$ is also known as the *0-dimensional sphere*, \mathbb{S}^0 . (For instance, it is the space of points at distance 1 from the origin in \mathbb{R}^1 , just as the topological 1-sphere is the space of points at distance 1 from the origin in \mathbb{R}^2 .) Thus, Lemma 6.5.1 can be phrased suggestively as $\Sigma \mathbb{S}^0 \simeq \mathbb{S}^1$. In fact, this pattern continues: we can define all the spheres inductively by

$$\mathbb{S}^0 := \mathbf{2} \quad \text{and} \quad \mathbb{S}^{n+1} := \Sigma \mathbb{S}^n. \quad (6.5.2)$$

We can even start one dimension lower by defining $\mathbb{S}^{-1} := \mathbf{0}$, and observe that $\Sigma \mathbf{0} \simeq \mathbf{2}$.

To prove carefully that this agrees with the definition of \mathbb{S}^n from the previous section would require making the latter more explicit. However, we can show that the recursive definition has the same universal property that we would expect the other one to have. If (A, a_0) and (B, b_0) are pointed types (with basepoints often left implicit), let $\text{Map}_*(A, B)$ denote the type of based maps:

$$\text{Map}_*(A, B) := \sum_{f:A \rightarrow B} (f(a_0) = b_0).$$

Note that any type A gives rise to a pointed type $A_+ := A + \mathbf{1}$ with basepoint $\text{inr}(\star)$; this is called *adjoining a disjoint basepoint*.

Lemma 6.5.3. *For a type A and a pointed type (B, b_0) , we have*

$$\text{Map}_*(A_+, B) \simeq (A \rightarrow B)$$

Note that on the right we have the ordinary type of *unbased* functions from A to B .

Proof. From left to right, given $f : A_+ \rightarrow B$ with $p : f(\text{inr}(\star)) = b_0$, we have $f \circ \text{inl} : A \rightarrow B$. And from right to left, given $g : A \rightarrow B$ we define $g' : A_+ \rightarrow B$ by $g'(\text{inl}(a)) := g(a)$ and $g'(\text{inr}(u)) := b_0$. We leave it to the reader to show that these are quasi-inverse operations. \square

In particular, note that $\mathbf{2} \simeq \mathbf{1}_+$. Thus, for any pointed type B we have

$$\text{Map}_*(\mathbf{2}, B) \simeq (\mathbf{1} \rightarrow B) \simeq B.$$

Now recall that the loop space operation Ω acts on pointed types, with definition $\Omega(A, a_0) := (a_0 =_A a_0, \text{refl}_{a_0})$. We can also make the suspension Σ act on pointed types, by $\Sigma(A, a_0) := (\Sigma A, \mathbf{N})$.

Lemma 6.5.4. *For pointed types (A, a_0) and (B, b_0) we have*

$$\text{Map}_*(\Sigma A, B) \simeq \text{Map}_*(A, \Omega B).$$

Proof. We first observe the following chain of equivalences:

$$\begin{aligned} \text{Map}_*(\Sigma A, B) &:= \sum_{f:\Sigma A \rightarrow B} (f(\mathbf{N}) = b_0) \\ &\simeq \sum_{f:\Sigma(b_n:B) \sum_{(b_s:B)} (A \rightarrow (b_n = b_s))} (\text{pr}_1(f) = b_0) \\ &\simeq \sum_{(b_n:B)} \sum_{(b_s:B)} (A \rightarrow (b_n = b_s)) \times (b_n = b_0) \\ &\simeq \sum_{(p:\sum_{(b_n:B)} (b_n = b_0))} \sum_{(b_s:B)} (A \rightarrow (\text{pr}_1(p) = b_s)) \\ &\simeq \sum_{b_s:B} (A \rightarrow (b_0 = b_s)) \end{aligned}$$

The first equivalence is by the universal property of suspensions, which says that

$$(\Sigma A \rightarrow B) \simeq \left(\sum_{(b_n:B)} \sum_{(b_s:B)} (A \rightarrow (b_n = b_s)) \right)$$

with the function from right to left given by the recursor (see Exercise 6.11). The second and third equivalences are by Exercise 2.10, along with a reordering of components. Finally, the last equivalence follows from Lemma 3.11.9, since by Lemma 3.11.8, $\sum_{(b_n:B)} (b_n = b_0)$ is contractible with center (b_0, refl_{b_0}) .

The proof is now completed by the following chain of equivalences:

$$\begin{aligned} \sum_{b_s:B} (A \rightarrow (b_0 = b_s)) &\simeq \sum_{(b_s:B)} \sum_{(g:A \rightarrow (b_0 = b_s))} \sum_{(q:b_0 = b_s)} (g(a_0) = q) \\ &\simeq \sum_{(r:\sum_{(b_s:B)} (b_0 = b_s))} \sum_{(g:A \rightarrow (b_0 = \text{pr}_1(r)))} (g(a_0) = \text{pr}_2(r)) \\ &\simeq \sum_{g:A \rightarrow (b_0 = b_0)} (g(a_0) = \text{refl}_{b_0}) \\ &\equiv \text{Map}_*(A, \Omega B). \end{aligned}$$

Similar to before, the first and last equivalences are by Lemmas 3.11.8 and 3.11.9, and the second is by Exercise 2.10 and reordering of components. \square

In particular, for the spheres defined as in (6.5.2) we have

$$\text{Map}_*(S^n, B) \simeq \text{Map}_*(S^{n-1}, \Omega B) \simeq \dots \simeq \text{Map}_*(\mathbf{2}, \Omega^n B) \simeq \Omega^n B.$$

Thus, these spheres S^n have the universal property that we would expect from the spheres defined directly in terms of n -fold loop spaces as in §6.4.

6.6 Cell complexes

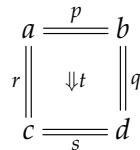
In classical topology, a *cell complex* is a space obtained by successively attaching discs along their boundaries. It is called a *CW complex* if the boundary of an n -dimensional disc is constrained to lie in the discs of dimension strictly less than n (the $(n-1)$ -skeleton).

Any finite CW complex can be presented as a higher inductive type, by turning n -dimensional discs into n -dimensional paths and partitioning the image of the attaching map into a source and a target, with each written as a composite of lower dimensional paths. Our explicit definitions of S^1 and S^2 in §6.4 had this form.

Another example is the torus T^2 , which is generated by:

- a point $b : T^2$,
- a path $p : b = b$,
- another path $q : b = b$, and
- a 2-path $t : p \cdot q = q \cdot p$.

Perhaps the easiest way to see that this is a torus is to start with a rectangle, having four corners a, b, c, d , four edges p, q, r, s , and an interior which is manifestly a 2-path t from $p \cdot q$ to $r \cdot s$:



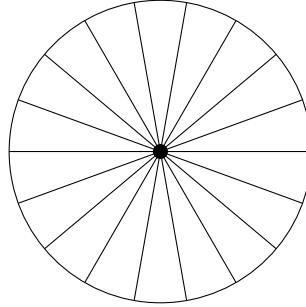


Figure 6.3: A 2-disc made out of a hub and spokes

Now identify the edge r with q and the edge s with p , resulting in also identifying all four corners. Topologically, this identification can be seen to produce a torus.

The induction principle for the torus is the trickiest of any we've written out so far. Given $P : T^2 \rightarrow \mathcal{U}$, for a section $\prod_{(x:T^2)} P(x)$ we require

- a point $b' : P(b)$,
- a path $p' : b' =_p^P b'$,
- a path $q' : b' =_q^P b'$, and
- a 2-path t' between the "composites" $p' \cdot q'$ and $q' \cdot p'$, lying over t .

In order to make sense of this last datum, we need a composition operation for dependent paths, but this is not hard to define. Then the induction principle gives a function $f : \prod_{(x:T^2)} P(x)$ such that $f(b) \equiv b'$ and $\text{apd}_f(p) = p'$ and $\text{apd}_f(q) = q'$ and something like " $\text{apd}_f^2(t) = t''$ ". However, this is not well-typed as it stands, firstly because the equalities $\text{apd}_f(p) = p'$ and $\text{apd}_f(q) = q'$ are not judgmental, and secondly because apd_f only preserves path concatenation up to homotopy. We leave the details to the reader (see Exercise 6.1).

Of course, another definition of the torus is $T^2 \equiv S^1 \times S^1$ (in Exercise 6.3 we ask the reader to verify the equivalence of the two). The cell-complex definition, however, generalizes easily to other spaces without such descriptions, such as the Klein bottle, the projective plane, etc. But it does get increasingly difficult to write down the induction principles, requiring us to define notions of dependent n -paths and of apd acting on n -paths. Fortunately, once we have the spheres in hand, there is a way around this.

6.7 Hubs and spokes

In topology, one usually speaks of building CW complexes by attaching n -dimensional discs along their $(n - 1)$ -dimensional boundary spheres. However, another way to express this is by gluing in the *cone* on an $(n - 1)$ -dimensional sphere. That is, we regard a disc as consisting of a cone point (or "hub"), with meridians (or "spokes") connecting that point to every point on the boundary, continuously, as shown in Figure 6.3.

We can use this idea to express higher inductive types containing n -dimensional path constructors for $n > 1$ in terms of ones containing only 1-dimensional path constructors. The point is that we can obtain an n -dimensional path as a continuous family of 1-dimensional paths parametrized by an $(n - 1)$ -dimensional object. The simplest $(n - 1)$ -dimensional object to use is the $(n - 1)$ -sphere, although in some cases a different one may be preferable. (Recall that

we were able to define the spheres in §6.5 inductively using suspensions, which involve only 1-dimensional path constructors. Indeed, suspension can also be regarded as an instance of this idea, since it involves a family of 1-dimensional paths parametrized by the type being suspended.)

For instance, the torus T^2 from the previous section could be defined instead to be generated by:

- a point $b : T^2$,
- a path $p : b = b$,
- another path $q : b = b$,
- a point $h : T^2$, and
- for each $x : \mathbb{S}^1$, a path $s(x) : f(x) = h$, where $f : \mathbb{S}^1 \rightarrow T^2$ is defined by $f(\text{base}) := b$ and $f(\text{loop}) := p \cdot q \cdot p^{-1} \cdot q^{-1}$.

The induction principle for this version of the torus says that given $P : T^2 \rightarrow \mathcal{U}$, for a section $\prod_{(x:T^2)} P(x)$ we require

- a point $b' : P(b)$,
- a path $p' : b' =_p^P b'$,
- a path $q' : b' =_q^P b'$,
- a point $h' : P(h)$, and
- for each $x : \mathbb{S}^1$, a path $g(x) =_{s(x)}^P h'$, where $g : \prod_{(x:\mathbb{S}^1)} P(f(x))$ is defined by $g(\text{base}) := b'$ and $\text{apd}_g(\text{loop}) := t(p' \cdot q' \cdot (p')^{-1} \cdot (q')^{-1})$. In the latter, \cdot denotes concatenation of dependent paths, and the definition of $t : (b' =_{f(\text{loop})}^P b') \simeq (b' =_{\text{loop}}^{P \circ f} b')$ is left to the reader.

Note that there is no need for dependent 2-paths or apd^2 . We leave it to the reader to write out the computation rules.

Remark 6.7.1. One might question the need for introducing the hub point h ; why couldn't we instead simply add paths continuously relating the boundary of the disc to a point *on* that boundary, as shown in Figure 6.4? However, this does not work without further modification. For if, given some $f : \mathbb{S}^1 \rightarrow X$, we give a path constructor connecting each $f(x)$ to $f(\text{base})$, then what we end up with is more like the picture in Figure 6.5 of a cone whose vertex is twisted around and glued to some point on its base. The problem is that the specified path from $f(\text{base})$ to itself may not be reflexivity. We could remedy the problem by adding a 2-dimensional path constructor to ensure this, but using a separate hub avoids the need for any path constructors of dimension above 1.

Remark 6.7.2. Note also that this “translation” of higher paths into 1-paths does not preserve judgmental computation rules for these paths, though it does preserve propositional ones.

6.8 Pushouts

From a category-theoretic point of view, one of the important aspects of any foundational system is the ability to construct limits and colimits. In set-theoretic foundations, these are limits and colimits of sets, whereas in our case they are limits and colimits of *types*. We have seen in §2.15

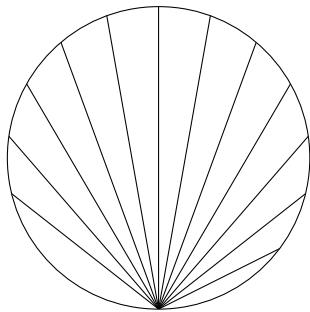


Figure 6.4: Hubless spokes

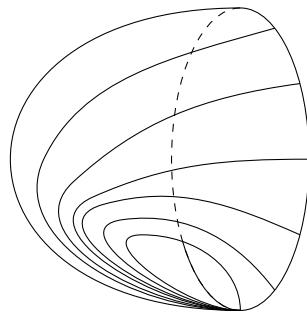


Figure 6.5: Hubless spokes, II

that cartesian product types have the correct universal property of a categorical product of types, and in Exercise 2.9 that coproduct types likewise have their expected universal property.

As remarked in §2.15, more general limits can be constructed using identity types and Σ -types, e.g. the pullback of $f : A \rightarrow C$ and $g : B \rightarrow C$ is $\sum_{(a:A)} \sum_{(b:B)} (f(a) = g(b))$ (see Exercise 2.11). However, more general *colimits* require identifying elements coming from different types, for which higher inductives are well-adapted. Since all our constructions are homotopy-invariant, all our colimits are necessarily *homotopy colimits*, but we drop the ubiquitous adjective in the interests of concision.

In this section we discuss *pushouts*, as perhaps the simplest and one of the most useful colimits. Indeed, one expects all finite colimits (for a suitable homotopical definition of “finite”) to be constructible from pushouts and finite coproducts. It is also possible to give a direct construction of more general colimits using higher inductive types, but this is somewhat technical, and also not completely satisfactory since we do not yet have a good fully general notion of homotopy coherent diagrams.

Suppose given a span of types and functions:

$$\mathcal{D} = \begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & & \\ A & & \end{array}$$

The **pushout** of this span is the higher inductive type $A \sqcup^C B$ presented by

- a function $\text{inl} : A \rightarrow A \sqcup^C B$,
- a function $\text{inr} : B \rightarrow A \sqcup^C B$, and
- for each $c : C$ a path $\text{glue}(c) : (\text{inl}(f(c)) = \text{inr}(g(c)))$.

In other words, $A \sqcup^C B$ is the disjoint union of A and B , together with for every $c : C$ a witness that $f(c)$ and $g(c)$ are equal. The recursion principle says that if D is another type, we can define a map $s : A \sqcup^C B \rightarrow D$ by defining

- for each $a : A$, the value of $s(\text{inl}(a)) : D$,
- for each $b : B$, the value of $s(\text{inr}(b)) : D$, and
- for each $c : C$, the value of $\text{ap}_s(\text{glue}(c)) : s(\text{inl}(f(c))) = s(\text{inr}(g(c)))$.

We leave it to the reader to formulate the induction principle. It also implies the uniqueness principle that if $s, s' : A \sqcup^C B \rightarrow D$ are two maps such that

$$\begin{aligned} s(\text{inl}(a)) &= s'(\text{inl}(a)) \\ s(\text{inr}(b)) &= s'(\text{inr}(b)) \\ \text{ap}_s(\text{glue}(c)) &= \text{ap}_{s'}(\text{glue}(c)) \quad (\text{modulo the previous two equalities}) \end{aligned}$$

for every a, b, c , then $s = s'$.

To formulate the universal property of a pushout, we introduce the following.

Definition 6.8.1. Given a span $\mathcal{D} = (A \xleftarrow{f} C \xrightarrow{g} B)$ and a type D , a **cocone under \mathcal{D} with vertex D** consists of functions $i : A \rightarrow D$ and $j : B \rightarrow D$ and a homotopy $h : \prod_{(c:C)} (i(f(c)) = j(g(c)))$:

$$\begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & h \not\equiv & \downarrow j \\ A & \xrightarrow{i} & D \end{array}$$

We denote by $\text{cocone}_{\mathcal{D}}(D)$ the type of all such cocones, i.e.

$$\text{cocone}_{\mathcal{D}}(D) := \sum_{(i:A \rightarrow D)} \sum_{(j:B \rightarrow D)} \prod_{(c:C)} (i(f(c)) = j(g(c))).$$

Of course, there is a canonical cocone under \mathcal{D} with vertex $A \sqcup^C B$ consisting of inl , inr , and glue .

$$\begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & \text{glue} \not\equiv & \downarrow \text{inr} \\ A & \xrightarrow{\text{inl}} & A \sqcup^C B \end{array}$$

The following lemma says that this is the universal such cocone.

Lemma 6.8.2. *For any type E , there is an equivalence*

$$(A \sqcup^C B \rightarrow E) \simeq \text{cocone}_{\mathcal{D}}(E).$$

Proof. Let's consider an arbitrary type $E : \mathcal{U}$. There is a canonical function c_{\sqcup} defined by

$$\begin{cases} (A \sqcup^C B \rightarrow E) & \longrightarrow \text{cocone}_{\mathcal{D}}(E) \\ t & \longmapsto (t \circ \text{inl}, t \circ \text{inr}, \text{ap}_t \circ \text{glue}) \end{cases}$$

We write informally $t \mapsto t \circ c_{\sqcup}$ for this function. We show that this is an equivalence.

Firstly, given a $c = (i, j, h) : \text{cocone}_{\mathcal{D}}(E)$, we need to construct a map $s(c)$ from $A \sqcup^C B$ to E .

$$\begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & h \not\equiv & \downarrow j \\ A & \xrightarrow{i} & E \end{array}$$

The map $s(c)$ is defined in the following way

$$\begin{aligned}s(c)(\text{inl}(a)) &:= i(a), \\ s(c)(\text{inr}(b)) &:= j(b), \\ \text{ap}_{s(c)}(\text{glue}(x)) &:= h(x).\end{aligned}$$

We have defined a map

$$\begin{cases} \text{cocone}_{\mathcal{D}}(E) & \longrightarrow (A \sqcup^C B \rightarrow E) \\ c & \longmapsto s(c) \end{cases}$$

and we need to prove that this map is an inverse to $t \mapsto t \circ c_{\sqcup}$. On the one hand, if $c = (i, j, h) : \text{cocone}_{\mathcal{D}}(E)$, we have

$$\begin{aligned}s(c) \circ c_{\sqcup} &= (s(c) \circ \text{inl}, s(c) \circ \text{inr}, \text{ap}_{s(c)} \circ \text{glue}) \\ &= (\lambda a. s(c)(\text{inl}(a)), \lambda b. s(c)(\text{inr}(b)), \lambda x. \text{ap}_{s(c)}(\text{glue}(x))) \\ &= (\lambda a. i(a), \lambda b. j(b), \lambda x. h(x)) \\ &\equiv (i, j, h) \\ &= c.\end{aligned}$$

On the other hand, if $t : A \sqcup^C B \rightarrow E$, we want to prove that $s(t \circ c_{\sqcup}) = t$. For $a : A$, we have

$$s(t \circ c_{\sqcup})(\text{inl}(a)) = t(\text{inl}(a))$$

because the first component of $t \circ c_{\sqcup}$ is $t \circ \text{inl}$. In the same way, for $b : B$ we have

$$s(t \circ c_{\sqcup})(\text{inr}(b)) = t(\text{inr}(b))$$

and for $x : C$ we have

$$\text{ap}_{s(t \circ c_{\sqcup})}(\text{glue}(x)) = \text{ap}_t(\text{glue}(x))$$

hence $s(t \circ c_{\sqcup}) = t$.

This proves that $c \mapsto s(c)$ is a quasi-inverse to $t \mapsto t \circ c_{\sqcup}$, as desired. \square

A number of standard homotopy-theoretic constructions can be expressed as (homotopy) pushouts.

- The pushout of the span $\mathbf{1} \leftarrow A \rightarrow \mathbf{1}$ is the **suspension** ΣA (see §6.5).
- The pushout of $A \xleftarrow{\text{pr}_1} A \times B \xrightarrow{\text{pr}_2} B$ is called the **join** of A and B , written $A * B$.
- The pushout of $\mathbf{1} \leftarrow A \xrightarrow{f} B$ is the **cone** or **cofiber** of f .
- If A and B are equipped with basepoints $a_0 : A$ and $b_0 : B$, then the pushout of $A \xleftarrow{a_0} \mathbf{1} \xrightarrow{b_0} B$ is the **wedge** $A \vee B$.
- If A and B are pointed as before, define $f : A \vee B \rightarrow A \times B$ by $f(\text{inl}(a)) := (a, b_0)$ and $f(\text{inr}(b)) := (a_0, b)$, with $f(\text{glue}) := \text{refl}_{(a_0, b_0)}$. Then the cone of f is called the **smash product** $A \wedge B$.

We will discuss pushouts further in Chapter 7 and ??.

Remark 6.8.3. As remarked in §3.7, the notations \wedge and \vee for the smash product and wedge of pointed spaces are also used in logic for “and” and “or”, respectively. Since types in homotopy type theory can behave either like spaces or like propositions, there is technically a potential for conflict — but since they rarely do both at once, context generally disambiguates. Furthermore, the smash product and wedge only apply to *pointed* spaces, while the only pointed mere proposition is $\top \equiv \mathbf{1}$ — and we have $\mathbf{1} \wedge \mathbf{1} = \mathbf{1}$ and $\mathbf{1} \vee \mathbf{1} = \mathbf{1}$ for either meaning of \wedge and \vee .

Remark 6.8.4. Note that colimits do not in general preserve truncatedness. For instance, S^0 and $\mathbf{1}$ are both sets, but the pushout of $\mathbf{1} \leftarrow S^0 \rightarrow \mathbf{1}$ is S^1 , which is not a set. If we are interested in colimits in the category of n -types, therefore (and, in particular, in the category of sets), we need to “truncate” the colimit somehow. We will return to this point in §6.9, Chapter 7, and ??.

6.9 Truncations

In §3.7 we introduced the propositional truncation as a new type forming operation; we now observe that it can be obtained as a special case of higher inductive types. This reduces the problem of understanding truncations to the problem of understanding higher inductives, which at least are amenable to a systematic treatment. It is also interesting because it provides our first example of a higher inductive type which is truly *recursive*, in that its constructors take inputs from the type being defined (as does the successor $\text{succ} : \mathbb{N} \rightarrow \mathbb{N}$).

Let A be a type; we define its propositional truncation $\|A\|$ to be the higher inductive type generated by:

- A function $|-| : A \rightarrow \|A\|$, and
- for each $x, y : \|A\|$, a path $x = y$.

Note that the second constructor is by definition the assertion that $\|A\|$ is a mere proposition. Thus, the definition of $\|A\|$ can be interpreted as saying that $\|A\|$ is freely generated by a function $A \rightarrow \|A\|$ and the fact that it is a mere proposition.

The recursion principle for this higher inductive definition is easy to write down: it says that given any type B together with

- a function $g : A \rightarrow B$, and
- for any $x, y : B$, a path $x =_B y$,

there exists a function $f : \|A\| \rightarrow B$ such that

- $f(|a|) \equiv g(a)$ for all $a : A$, and
- for any $x, y : \|A\|$, the function ap_f takes the specified path $x = y$ in $\|A\|$ to the specified path $f(x) = f(y)$ in B (propositionally).

These are exactly the hypotheses that we stated in §3.7 for the recursion principle of propositional truncation — a function $A \rightarrow B$ such that B is a mere proposition — and the first part of the conclusion is exactly what we stated there as well. The second part (the action of ap_f) was not mentioned previously, but it turns out to be vacuous in this case, because B is a mere proposition, so *any* two paths in it are automatically equal.

There is also an induction principle for $\|A\|$, which says that given any $B : \|A\| \rightarrow \mathcal{U}$ together with

- a function $g : \prod_{(a:A)} B(|a|)$, and

- for any $x, y : \|A\|$ and $u : B(x)$ and $v : B(y)$, a dependent path $q : u =_{p(x,y)}^B v$, where $p(x, y)$ is the path coming from the second constructor of $\|A\|$,

there exists $f : \prod_{(x:\|A\|)} B(x)$ such that $f(|a|) \equiv g(a)$ for $a : A$, and also another computation rule. However, because there can be at most one function between any two mere propositions (up to homotopy), this induction principle is not really useful (see also Exercise 3.17).

We can, however, extend this idea to construct similar truncations landing in n -types, for any n . For instance, we might define the 0 -truncation $\|A\|_0$ to be generated by

- A function $| - |_0 : A \rightarrow \|A\|_0$, and
- For each $x, y : \|A\|_0$ and each $p, q : x = y$, a path $p = q$.

Then $\|A\|_0$ would be freely generated by a function $A \rightarrow \|A\|_0$ together with the assertion that $\|A\|_0$ is a set. A natural induction principle for it would say that given $B : \|A\|_0 \rightarrow \mathcal{U}$ together with

- a function $g : \prod_{(a:A)} B(|a|_0)$, and
- for any $x, y : \|A\|_0$ with $z : B(x)$ and $w : B(y)$, and each $p, q : x = y$ with $r : z =_p^B w$ and $s : z =_q^B w$, a 2-path $v : r =_{u(x,y,p,q)}^{z =_w^B} s$, where $u(x, y, p, q) : p = q$ is obtained from the second constructor of $\|A\|_0$,

there exists $f : \prod_{(x:\|A\|_0)} B(x)$ such that $f(|a|_0) \equiv g(a)$ for all $a : A$, and also $\text{apd}_f^2(u(x, y, p, q))$ is the 2-path specified above. (As in the propositional case, the latter condition turns out to be uninteresting.) From this, however, we can prove a more useful induction principle.

Lemma 6.9.1. *Suppose given $B : \|A\|_0 \rightarrow \mathcal{U}$ together with $g : \prod_{(a:A)} B(|a|_0)$, and assume that each $B(x)$ is a set. Then there exists $f : \prod_{(x:\|A\|_0)} B(x)$ such that $f(|a|_0) \equiv g(a)$ for all $a : A$.*

Proof. It suffices to construct, for any x, y, z, w, p, q, r, s as above, a 2-path $v : r =_{u(x,y,p,q)}^{z =_w^B} s$. However, by the definition of dependent 2-paths, this is an ordinary 2-path in the fiber $B(y)$. Since $B(y)$ is a set, a 2-path exists between any two parallel paths. \square

This implies the expected universal property.

Lemma 6.9.2. *For any set B and any type A , composition with $| - |_0 : A \rightarrow \|A\|_0$ determines an equivalence*

$$(\|A\|_0 \rightarrow B) \simeq (A \rightarrow B).$$

Proof. The special case of Lemma 6.9.1 when B is the constant family gives a map from right to left, which is a right inverse to the “compose with $| - |_0$ ” function from left to right. To show that it is also a left inverse, let $h : \|A\|_0 \rightarrow B$, and define $h' : \|A\|_0 \rightarrow B$ by applying Lemma 6.9.1 to the composite $a \mapsto h(|a|_0)$. Thus, $h'(|a|_0) = h(|a|_0)$.

However, since B is a set, for any $x : \|A\|_0$ the type $h(x) = h'(x)$ is a mere proposition, and hence also a set. Therefore, by Lemma 6.9.1, the observation that $h'(|a|_0) = h(|a|_0)$ for any $a : A$ implies $h(x) = h'(x)$ for any $x : \|A\|_0$, and hence $h = h'$. \square

For instance, this enables us to construct colimits of sets. We have seen that if $A \xleftarrow{f} C \xrightarrow{g} B$ is a span of sets, then the pushout $A \sqcup^C B$ may no longer be a set. (For instance, if A and B are $\mathbf{1}$ and C is $\mathbf{2}$, then the pushout is S^1 .) However, we can construct a pushout that is a set, and has the expected universal property with respect to other sets, by truncating.

Lemma 6.9.3. Let $A \xleftarrow{f} C \xrightarrow{g} B$ be a span of sets. Then for any set E , there is a canonical equivalence

$$\left(\|A \sqcup^C B\|_0 \rightarrow E \right) \simeq \text{cocone}_{\mathcal{D}}(E).$$

Proof. Compose the equivalences in Lemmas 6.8.2 and 6.9.2. \square

We refer to $\|A \sqcup^C B\|_0$ as the **set-pushout** of f and g , to distinguish it from the (homotopy) pushout $A \sqcup^C B$. Alternatively, we could modify the definition of the pushout in §6.8 to include the 0-truncation constructor directly, avoiding the need to truncate afterwards. Similar remarks apply to any sort of colimit of sets; we will explore this further in ??.

However, while the above definition of the 0-truncation works — it gives what we want, and is consistent — it has a couple of issues. Firstly, it doesn't fit so nicely into the general theory of higher inductive types. In general, it is tricky to deal directly with constructors such as the second one we have given for $\|A\|_0$, whose *inputs* involve not only elements of the type being defined, but paths in it.

This can be gotten round fairly easily, however. Recall in §5.1 we mentioned that we can allow a constructor of an inductive type W to take “infinitely many arguments” of type W by having it take a single argument of type $\mathbb{N} \rightarrow W$. There is a general principle behind this: to model a constructor with funny-looking inputs, use an auxiliary inductive type (such as \mathbb{N}) to parametrize them, reducing the input to a simple function with inductive domain.

For the 0-truncation, we can consider the auxiliary *higher* inductive type S generated by two points $a, b : S$ and two paths $p, q : a = b$. Then the fishy-looking constructor of $\|A\|_0$ can be replaced by the unobjectionable

- For every $f : S \rightarrow \|A\|_0$, a path $\text{ap}_f(p) = \text{ap}_f(q)$.

Since to give a map out of S is the same as to give two points and two parallel paths between them, this yields the same induction principle.

A more serious problem with our current definition of 0-truncation, however, is that it doesn't generalize very well. If we want to describe a notion of definition of “ n -truncation” into n -types uniformly for all $n : \mathbb{N}$, then this approach is unfeasible, since the second constructor would need a number of arguments that increases with n . In §7.3, therefore, we will use a different idea to construct these, based on the observation that the type S introduced above is equivalent to the circle S^1 . This includes the 0-truncation as a special case, and satisfies generalized versions of Lemmas 6.9.1 and 6.9.2.

6.10 Quotients

A particularly important sort of colimit of sets is the *quotient* by a relation. That is, let A be a set and $R : A \times A \rightarrow \text{Prop}$ a family of mere propositions (a **mere relation**). Its quotient should be the set-coequalizer of the two projections

$$\sum_{(a,b:A)} R(a,b) \rightrightarrows A.$$

We can also describe this directly, as the higher inductive type A/R generated by

- A function $q : A \rightarrow A/R$;
- For each $a, b : A$ such that $R(a, b)$, an equality $q(a) = q(b)$; and

- The 0-truncation constructor: for all $x, y : A/R$ and $r, s : x = y$, we have $r = s$.

We will sometimes refer to this higher inductive type A/R as the **set-quotient** of A by R , to emphasize that it produces a set by definition. (There are more general notions of “quotient” in homotopy theory, but they are mostly beyond the scope of this book. However, in ?? we will consider the “quotient” of a type by a 1-groupoid, which is the next level up from set-quotients.)

Remark 6.10.1. It is not actually necessary for the definition of set-quotients, and most of their properties, that A be a set. However, this is generally the case of most interest.

Lemma 6.10.2. *The function $q : A \rightarrow A/R$ is surjective.*

Proof. We must show that for any $x : A/R$ there merely exists an $a : A$ with $q(a) = x$. We use the induction principle of A/R . The first case is trivial: if x is $q(a)$, then of course there merely exists an a such that $q(a) = q(a)$. And since the goal is a mere proposition, it automatically respects all path constructors, so we are done. \square

We can now prove that the set-quotient has the expected universal property of a (set-)coequalizer.

Lemma 6.10.3. *For any set B , precomposing with q yields an equivalence*

$$(A/R \rightarrow B) \simeq \left(\sum_{(f:A \rightarrow B)} \prod_{(a,b:A)} R(a,b) \rightarrow (f(a) = f(b)) \right).$$

Proof. The quasi-inverse of $- \circ q$, going from right to left, is just the recursion principle for A/R . That is, given $f : A \rightarrow B$ such that $\prod_{(a,b:A)} R(a,b) \rightarrow (f(a) = f(b))$, we define $\bar{f} : A/R \rightarrow B$ by $\bar{f}(q(a)) := f(a)$. This defining equation says precisely that $(f \mapsto \bar{f})$ is a right inverse to $(- \circ q)$.

For it to also be a left inverse, we must show that for any $g : A/R \rightarrow B$ and $x : A/R$ we have $g(x) = \overline{g \circ q}(x)$. However, by Lemma 6.10.2 there merely exists a such that $q(a) = x$. Since our desired equality is a mere proposition, we may assume there purely exists such an a , in which case $g(x) = g(q(a)) = \overline{g \circ q}(q(a)) = \overline{g \circ q}(x)$. \square

Of course, classically the usual case to consider is when R is an **equivalence relation**, i.e. we have

- **reflexivity:** $\prod_{(a:A)} R(a,a)$,
- **symmetry:** $\prod_{(a,b:A)} R(a,b) \rightarrow R(b,a)$, and
- **transitivity:** $\prod_{(a,b,c:C)} R(a,b) \times R(b,c) \rightarrow R(a,c)$.

In this case, the set-quotient A/R has additional good properties, as we will see in ??: for instance, we have $R(a,b) \simeq (q(a) =_{A/R} q(b))$. We often write an equivalence relation $R(a,b)$ infix as $a \sim b$.

The quotient by an equivalence relation can also be constructed in other ways. The set theoretic approach is to consider the set of equivalence classes, as a subset of the power set of A . We can mimic this “impredicative” construction in type theory as well.

Definition 6.10.4. A predicate $P : A \rightarrow \text{Prop}$ is an **equivalence class** of a relation $R : A \times A \rightarrow \text{Prop}$ if there merely exists an $a : A$ such that for all $b : A$ we have $R(a,b) \simeq P(b)$.

As R and P are mere propositions, the equivalence $R(a,b) \simeq P(b)$ is the same thing as implications $R(a,b) \rightarrow P(b)$ and $P(b) \rightarrow R(a,b)$. And of course, for any $a : A$ we have the canonical equivalence class $P_a(b) := R(a,b)$.

Definition 6.10.5. We define

$$A // R := \{ P : A \rightarrow \text{Prop} \mid P \text{ is an equivalence class of } R \}.$$

The function $q' : A \rightarrow A // R$ is defined by $q'(a) := P_a$.

Theorem 6.10.6. For any equivalence relation R on A , the type $A // R$ is equivalent to the set-quotient A/R .

Proof. First, note that if $R(a, b)$, then since R is an equivalence relation we have $R(a, c) \Leftrightarrow R(b, c)$ for any $c : A$. Thus, $R(a, c) = R(b, c)$ by univalence, hence $P_a = P_b$ by function extensionality, i.e. $q'(a) = q'(b)$. Therefore, by Lemma 6.10.3 we have an induced map $f : A/R \rightarrow A // R$ such that $f \circ q = q'$.

We show that f is injective and surjective, hence an equivalence. Surjectivity follows immediately from the fact that q' is surjective, which in turn is true essentially by definition of $A // R$. For injectivity, if $f(x) = f(y)$, then to show the mere proposition $x = y$, by surjectivity of q we may assume $x = q(a)$ and $y = q(b)$ for some $a, b : A$. Then $R(a, c) = f(q(a))(c) = f(q(b))(c) = R(b, c)$ for any $c : A$, and in particular $R(a, b) = R(b, b)$. But $R(b, b)$ is inhabited, since R is an equivalence relation, hence so is $R(a, b)$. Thus $q(a) = q(b)$ and so $x = y$. \square

In ?? we will give an alternative proof of this theorem. Note that unlike A/R , the construction $A // R$ raises universe level: if $A : \mathcal{U}_i$ and $R : A \rightarrow A \rightarrow \text{Prop}_{\mathcal{U}_i}$, then in the definition of $A // R$ we must also use $\text{Prop}_{\mathcal{U}_i}$ to include all the equivalence classes, so that $A // R : \mathcal{U}_{i+1}$. Of course, we can avoid this if we assume the propositional resizing axiom from §3.5.

Remark 6.10.7. The previous two constructions provide quotients in generality, but in particular cases there may be easier constructions. For instance, we may define the integers \mathbb{Z} as a set-quotient

$$\mathbb{Z} := (\mathbb{N} \times \mathbb{N}) / \sim$$

where \sim is the equivalence relation defined by

$$(a, b) \sim (c, d) := (a + d = b + c).$$

In other words, a pair (a, b) represents the integer $a - b$. In this case, however, there are *canonical representatives* of the equivalence classes: those of the form $(n, 0)$ or $(0, n)$.

The following lemma says that when this sort of thing happens, we don't need either general construction of quotients. (A function $r : A \rightarrow A$ is called **idempotent** if $r \circ r = r$.)

Lemma 6.10.8. Suppose \sim is a relation on a set A , and there exists an idempotent $r : A \rightarrow A$ such that $(r(x) = r(y)) \simeq (x \sim y)$ for all $x, y : A$. (This implies \sim is an equivalence relation.) Then the type

$$(A / \sim) := \left(\sum_{x:A} r(x) = x \right)$$

satisfies the universal property of the set-quotient of A by \sim , and hence is equivalent to it. In other words, there is a map $q : A \rightarrow (A / \sim)$ such that for every set B , precomposition with q induces an equivalence

$$\left((A / \sim) \rightarrow B \right) \simeq \left(\sum_{(g:A \rightarrow B)} \prod_{(x,y:A)} (x \sim y) \rightarrow (g(x) = g(y)) \right). \quad (6.10.9)$$

Proof. Let $i : \prod_{(x:A)} r(r(x)) = r(x)$ witness idempotence of r . The map $q : A \rightarrow (A/\sim)$ is defined by $q(x) := (r(x), i(x))$. Note that since A is a set, we have $q(x) = q(y)$ if and only if $r(x) = r(y)$, hence (by assumption) if and only if $x \sim y$. We define a map e from left to right in (6.10.9) by

$$e(f) := (f \circ q, _),$$

where the underscore $_$ denotes the following proof: if $x, y : A$ and $x \sim y$, then $q(x) = q(y)$ as observed above, hence $f(q(x)) = f(q(y))$. To see that e is an equivalence, consider the map e' in the opposite direction defined by

$$e'(g, s)(x, p) := g(x).$$

Given any $f : (A/\sim) \rightarrow B$,

$$e'(e(f))(x, p) \equiv f(q(x)) \equiv f(r(x), i(x)) = f(x, p)$$

where the last equality holds because $p : r(x) = x$ and so $(x, p) = (r(x), i(x))$ because A is a set. Similarly we compute

$$e(e'(g, s)) \equiv e(g \circ \text{pr}_1) \equiv (g \circ \text{pr}_1 \circ q, _).$$

Because B is a set we need not worry about the $_$ part, while for the first component we have

$$g(\text{pr}_1(q(x))) \equiv g(r(x)) = g(x),$$

where the last equation holds because $r(x) \sim x$, and g respects \sim by the assumption s . \square

Corollary 6.10.10. *Suppose $p : A \rightarrow B$ is a retraction between sets. Then B is the quotient of A by the equivalence relation \sim defined by*

$$(a_1 \sim a_2) := (p(a_1) = p(a_2)).$$

Proof. Suppose $s : B \rightarrow A$ is a section of p . Then $s \circ p : A \rightarrow A$ is an idempotent which satisfies the condition of Lemma 6.10.8 for this \sim , and s induces an isomorphism from B to its set of fixed points. \square

Remark 6.10.11. Lemma 6.10.8 applies to \mathbb{Z} with the idempotent $r : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \times \mathbb{N}$ defined by

$$r(a, b) = \begin{cases} (a - b, 0) & \text{if } a \geq b, \\ (0, b - a) & \text{otherwise.} \end{cases}$$

(This is a valid definition even constructively, since the relation \geq on \mathbb{N} is decidable.) Thus a non-negative integer is canonically represented as $(k, 0)$ and a non-positive one by $(0, m)$, for $k, m : \mathbb{N}$. This division into cases implies the following “induction principle” for integers, which will be useful in ???. (As usual, we identify a natural number n with the corresponding non-negative integer, i.e. with the image of $(n, 0) : \mathbb{N} \times \mathbb{N}$ in \mathbb{Z} .)

Lemma 6.10.12. *Suppose $P : \mathbb{Z} \rightarrow \mathcal{U}$ is a type family and that we have*

- $d_0 : P(0)$,
- $d_+ : \prod_{(n:\mathbb{N})} P(n) \rightarrow P(\text{succ}(n))$, and
- $d_- : \prod_{(n:\mathbb{N})} P(-n) \rightarrow P(-\text{succ}(n))$.

Then we have $f : \prod_{(z:\mathbb{Z})} P(z)$ such that

- $f(0) = d_0$,
- $f(\text{succ}(n)) = d_+(n, f(n))$ for all $n : \mathbb{N}$, and
- $f(-\text{succ}(n)) = d_-(n, f(-n))$ for all $n : \mathbb{N}$.

Proof. For purposes of this proof, let \mathbb{Z} denote $\sum_{(x:\mathbb{N} \times \mathbb{N})} (r(x) = x)$, where r is the above idempotent. (We can then transport the result to any equivalent definition of \mathbb{Z} .) Let $q : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{Z}$ be the quotient map, defined by $q(x) = (r(x), i(x))$ as in Lemma 6.10.8. Now define $Q := P \circ q : \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{U}$. By transporting the given data across appropriate equalities, we obtain

$$\begin{aligned} d'_0 &: Q(0, 0) \\ d'_+ &: \prod_{n:\mathbb{N}} Q(n, 0) \rightarrow Q(\text{succ}(n), 0) \\ d'_- &: \prod_{n:\mathbb{N}} Q(0, n) \rightarrow Q(0, \text{succ}(n)). \end{aligned}$$

Note also that since $q(n, m) = q(\text{succ}(n), \text{succ}(m))$, we have an induced equivalence

$$e_{n,m} : Q(n, m) \simeq Q(\text{succ}(n), \text{succ}(m)).$$

We can then construct $g : \prod_{(x:\mathbb{N} \times \mathbb{N})} Q(x)$ by double induction on x :

$$\begin{aligned} g(0, 0) &:= d'_0, \\ g(\text{succ}(n), 0) &:= d'_+(n, g(n, 0)), \\ g(0, \text{succ}(m)) &:= d'_-(m, g(0, m)), \\ g(\text{succ}(n), \text{succ}(m)) &:= e_{n,m}(g(n, m)). \end{aligned}$$

Now we have $\text{pr}_1 : \mathbb{Z} \rightarrow \mathbb{N} \times \mathbb{N}$, with the property that $q \circ \text{pr}_1 = \text{id}$. In particular, therefore, we have $Q \circ \text{pr}_1 = P$, and hence a family of equivalences $s : \prod_{(z:\mathbb{Z})} Q(\text{pr}_1(z)) \simeq P(z)$. Thus, we can define $f(z) = s(z, g(\text{pr}_1(z)))$ to obtain $f : \prod_{(z:\mathbb{Z})} P(z)$, and verify the desired equalities. \square

We will sometimes denote a function $f : \prod_{(z:\mathbb{Z})} P(z)$ obtained from Lemma 6.10.12 with a pattern-matching syntax, involving the three cases d_0 , d_+ , and d_- :

$$\begin{aligned} f(0) &:= d_0 \\ f(\text{succ}(n)) &:= d_+(n, f(n)) \\ f(-\text{succ}(n)) &:= d_-(n, f(-n)) \end{aligned}$$

We use $:=$ rather than \equiv , as we did for the path constructors of higher inductive types, to indicate that the “computation” rules implied by Lemma 6.10.12 are only propositional equalities. For example, in this way we can define the n -fold concatenation of a loop for any integer n .

Corollary 6.10.13. *Let A be a type with $a : A$ and $p : a = a$. There is a function $\prod_{(n:\mathbb{Z})} (a = a)$, denoted $n \mapsto p^n$, defined by*

$$\begin{aligned} p^0 &:= \text{refl}_a \\ p^{n+1} &:= p^n \cdot p && \text{for } n \geq 0 \\ p^{n-1} &:= p^n \cdot p^{-1} && \text{for } n \leq 0. \end{aligned}$$

We will discuss the integers further in §6.11 and ??.

6.11 Algebra

In addition to constructing higher-dimensional objects such as spheres and cell complexes, higher inductive types are also very useful even when working only with sets. We have seen one example already in Lemma 6.9.3: they allow us to construct the colimit of any diagram of sets, which is not possible in the base type theory of Chapter 1. Higher inductive types are also very useful when we study sets with algebraic structure.

As a running example in this section, we consider *groups*, which are familiar to most mathematicians and exhibit the essential phenomena (and will be needed in later chapters). However, most of what we say applies equally well to any sort of algebraic structure.

Definition 6.11.1. A **monoid** is a set G together with

- a *multiplication* function $G \times G \rightarrow G$, written infix as $(x, y) \mapsto x \cdot y$; and
- a *unit* element $e : G$; such that
- for any $x : G$, we have $x \cdot e = x$ and $e \cdot x = x$; and
- for any $x, y, z : G$, we have $x \cdot (y \cdot z) = (x \cdot y) \cdot z$.

A **group** is a monoid G together with

- an *inversion* function $i : G \rightarrow G$, written $x \mapsto x^{-1}$; such that
- for any $x : G$ we have $x \cdot x^{-1} = e$ and $x^{-1} \cdot x = e$.

Remark 6.11.2. Note that we require a group to be a set. We could consider a more general notion of “ ∞ -group” which is not a set, but this would take us further afield than is appropriate at the moment. With our current definition, we may expect the resulting “group theory” to behave similarly to the way it does in set-theoretic mathematics (with the caveat that, unless we assume LEM, it will be “constructive” group theory).

Example 6.11.3. The natural numbers \mathbb{N} are a monoid under addition, with unit 0, and also under multiplication, with unit 1. If we define the arithmetical operations on the integers \mathbb{Z} in the obvious way, then as usual they are a group under addition and a monoid under multiplication (and, of course, a ring). For instance, if $u, v \in \mathbb{Z}$ are represented by (a, b) and (c, d) , respectively, then $u + v$ is represented by $(a + c, b + d)$, $-u$ is represented by (b, a) , and uv is represented by $(ac + bd, ad + bc)$.

Example 6.11.4. We essentially observed in §2.1 that if (A, a) is a pointed type, then its loop space $\Omega(A, a) := (a =_A a)$ has all the structure of a group, except that it is not in general a set. It should be an “ ∞ -group” in the sense mentioned in Remark 6.11.2, but we can also make it a group by truncation. Specifically, we define the **fundamental group** of A based at $a : A$ to be

$$\pi_1(A, a) := \|\Omega(A, a)\|_0.$$

This inherits a group structure; for instance, the multiplication $\pi_1(A, a) \times \pi_1(A, a) \rightarrow \pi_1(A, a)$ is defined by double induction on truncation from the concatenation of paths.

More generally, the n^{th} **homotopy group** of (A, a) is $\pi_n(A, a) := \|\Omega^n(A, a)\|_0$. Then $\pi_n(A, a) = \pi_1(\Omega^{n-1}(A, a))$ for $n \geq 1$, so it is also a group. (When $n = 0$, we have $\pi_0(A) \equiv \|A\|_0$, which is not a group.) Moreover, the Eckmann–Hilton argument (Theorem 2.1.6) implies that if $n \geq 2$, then $\pi_n(A, a)$ is an *abelian* group, i.e. we have $x \cdot y = y \cdot x$ for all x, y . ?? will be largely the study of these groups.

One important notion in group theory is that of the *free group* generated by a set, or more generally of a group *presented* by generators and relations. It is well-known in type theory that *some* free algebraic objects can be defined using *ordinary* inductive types. For instance, the free monoid on a set A can be identified with the type $\text{List}(A)$ of *finite lists* of elements of A , which is inductively generated by

- a constructor $\text{nil} : \text{List}(A)$, and
- for each $\ell : \text{List}(A)$ and $a : A$, an element $\text{cons}(a, \ell) : \text{List}(A)$.

We have an obvious inclusion $\eta : A \rightarrow \text{List}(A)$ defined by $a \mapsto \text{cons}(a, \text{nil})$. The monoid operation on $\text{List}(A)$ is concatenation, defined recursively by

$$\begin{aligned}\text{nil} \cdot \ell &:= \ell \\ \text{cons}(a, \ell_1) \cdot \ell_2 &:= \text{cons}(a, \ell_1 \cdot \ell_2).\end{aligned}$$

It is straightforward to prove, using the induction principle for $\text{List}(A)$, that $\text{List}(A)$ is a set and that concatenation of lists is associative and has nil as a unit. Thus, $\text{List}(A)$ is a monoid.

Lemma 6.11.5. *For any set A , the type $\text{List}(A)$ is the free monoid on A . In other words, for any monoid G , composition with η is an equivalence*

$$\hom_{\text{Monoid}}(\text{List}(A), G) \simeq (A \rightarrow G),$$

where $\hom_{\text{Monoid}}(-, -)$ denotes the set of monoid homomorphisms (functions which preserve the multiplication and unit).

Proof. Given $f : A \rightarrow G$, we define $\bar{f} : \text{List}(A) \rightarrow G$ by recursion:

$$\begin{aligned}\bar{f}(\text{nil}) &:= e \\ \bar{f}(\text{cons}(a, \ell)) &:= f(a) \cdot \bar{f}(\ell).\end{aligned}$$

It is straightforward to prove by induction that \bar{f} is a monoid homomorphism, and that $f \mapsto \bar{f}$ is a quasi-inverse of $(-\circ\eta)$; see Exercise 6.8. \square

This construction of the free monoid is possible essentially because elements of the free monoid have computable canonical forms (namely, finite lists). However, elements of other free (and presented) algebraic structures — such as groups — do not in general have *computable* canonical forms. For instance, equality of words in group presentations is algorithmically undecidable. However, we can still describe free algebraic objects as *higher* inductive types, by simply asserting all the axiomatic equations as path constructors.

For example, let A be a set, and define a higher inductive type $F(A)$ with the following generators.

- A function $\eta : A \rightarrow F(A)$.
- A function $m : F(A) \times F(A) \rightarrow F(A)$.
- An element $e : F(A)$.
- A function $i : F(A) \rightarrow F(A)$.
- For each $x, y, z : F(A)$, an equality $m(x, m(y, z)) = m(m(x, y), z)$.
- For each $x : F(A)$, equalities $m(x, e) = x$ and $m(e, x) = x$.

- For each $x : F(A)$, equalities $m(x, i(x)) = e$ and $m(i(x), x) = e$.
- The 0-truncation constructor: for any $x, y : F(A)$ and $p, q : x = y$, we have $p = q$.

The first constructor says that A maps to $F(A)$. The next three give $F(A)$ the operations of a group: multiplication, an identity element, and inversion. The three constructors after that assert the axioms of a group: associativity, unitality, and inverses. Finally, the last constructor asserts that $F(A)$ is a set.

Therefore, $F(A)$ is a group. It is also straightforward to prove:

Theorem 6.11.6. *$F(A)$ is the free group on A . In other words, for any (set) group G , composition with $\eta : A \rightarrow F(A)$ determines an equivalence*

$$\text{hom}_{\text{Group}}(F(A), G) \simeq (A \rightarrow G)$$

where $\text{hom}_{\text{Group}}(-, -)$ denotes the set of group homomorphisms between two groups.

Proof. The recursion principle of the higher inductive type $F(A)$ says precisely that if G is a group and we have $f : A \rightarrow G$, then we have $\bar{f} : F(A) \rightarrow G$. Its computation rules say that $\bar{f} \circ \eta \equiv f$, and that \bar{f} is a group homomorphism. Thus, $(-\circ\eta) : \text{hom}_{\text{Group}}(F(A), G) \rightarrow (A \rightarrow G)$ has a right inverse. It is straightforward to use the induction principle of $F(A)$ to show that this is also a left inverse. \square

It is worth taking a step back to consider what we have just done. We have proven that the free group on any set exists *without* giving an explicit construction of it. Essentially all we had to do was write down the universal property that it should satisfy. In set theory, we could achieve a similar result by appealing to black boxes such as the adjoint functor theorem; type theory builds such constructions into the foundations of mathematics.

Of course, it is sometimes also useful to have a concrete description of free algebraic structures. In the case of free groups, we can provide one, using quotients. Consider $\text{List}(A + A)$, where in $A + A$ we write $\text{inl}(a)$ as a , and $\text{inr}(a)$ as \hat{a} (intended to stand for the formal inverse of a). The elements of $\text{List}(A + A)$ are *words* for the free group on A .

Theorem 6.11.7. *Let A be a set, and let $F'(A)$ be the set-quotient of $\text{List}(A + A)$ by the following relations.*

$$\begin{aligned} (\dots, a_1, a_2, \hat{a}_2, a_3, \dots) &= (\dots, a_1, a_3, \dots) \\ (\dots, a_1, \hat{a}_2, a_2, a_3, \dots) &= (\dots, a_1, a_3, \dots). \end{aligned}$$

Then $F'(A)$ is also the free group on the set A .

Proof. First we show that $F'(A)$ is a group. We have seen that $\text{List}(A + A)$ is a monoid; we claim that the monoid structure descends to the quotient. We define $F'(A) \times F'(A) \rightarrow F'(A)$ by double quotient recursion; it suffices to check that the equivalence relation generated by the given relations is preserved by concatenation of lists. Similarly, we prove the associativity and unit laws by quotient induction.

In order to define inverses in $F'(A)$, we first define $\text{reverse} : \text{List}(B) \rightarrow \text{List}(B)$ by recursion on lists:

$$\begin{aligned} \text{reverse}(\text{nil}) &\coloneqq \text{nil}, \\ \text{reverse}(\text{cons}(b, \ell)) &\coloneqq \text{reverse}(\ell) \cdot \text{cons}(b, \text{nil}). \end{aligned}$$

Now we define $i : F'(A) \rightarrow F'(A)$ by quotient recursion, acting on a list $\ell : \text{List}(A + A)$ by switching the two copies of A and reversing the list. This preserves the relations, hence descends to the quotient. And we can prove that $i(x) \cdot x = e$ for $x : F'(A)$ by induction. First, quotient induction allows us to assume x comes from $\ell : \text{List}(A + A)$, and then we can do list induction; if we write $q : \text{List}(A + A) \rightarrow F'(A)$ for the quotient map, the cases are

$$\begin{aligned} i(q(\text{nil})) \cdot q(\text{nil}) &= q(\text{nil}) \cdot q(\text{nil}) \\ &= q(\text{nil}) \\ i(q(\text{cons}(a, \ell))) \cdot q(\text{cons}(a, \ell)) &= i(q(\ell)) \cdot q(\text{cons}(\hat{a}, \text{nil})) \cdot q(\text{cons}(a, \ell)) \\ &= i(q(\ell)) \cdot q(\text{cons}(\hat{a}, \text{cons}(a, \ell))) \\ &= i(q(\ell)) \cdot q(\ell) \\ &= q(\text{nil}). \end{aligned} \quad (\text{by the inductive hypothesis})$$

(We have omitted a number of fairly evident lemmas about the behavior of concatenation of lists, etc.)

This completes the proof that $F'(A)$ is a group. Now if G is any group with a function $f : A \rightarrow G$, we can define $A + A \rightarrow G$ to be f on the first copy of A and f composed with the inversion map of G on the second copy. Now the fact that G is a monoid yields a monoid homomorphism $\text{List}(A + A) \rightarrow G$. And since G is a group, this map respects the relations, hence descends to a map $F'(A) \rightarrow G$. It is straightforward to prove that this is a group homomorphism, and the unique one which restricts to f on A . \square

If A has decidable equality (such as if we assume excluded middle), then the quotient defining $F'(A)$ can be obtained from an idempotent as in Lemma 6.10.8. We define a word, which we recall is just an element of $\text{List}(A + A)$, to be **reduced** if it contains no adjacent pairs of the form (a, \hat{a}) or (\hat{a}, a) . When A has decidable equality, it is straightforward to define the **reduction** of a word, which is an idempotent generating the appropriate quotient; we leave the details to the reader.

If $A := \mathbf{1}$, which has decidable equality, a reduced word must consist either entirely of $*$'s or entirely of $\hat{*}$'s. Thus, the free group on $\mathbf{1}$ is equivalent to the integers \mathbb{Z} , with 0 corresponding to nil, the positive integer n corresponding to a reduced word of $n *$'s, and the negative integer $(-n)$ corresponding to a reduced word of $n \hat{*}$'s. One could also, of course, show directly that \mathbb{Z} has the universal property of $F(\mathbf{1})$.

Remark 6.11.8. Nowhere in the construction of $F(A)$ and $F'(A)$, and the proof of their universal properties, did we use the assumption that A is a set. Thus, we can actually construct the free group on an arbitrary type. Comparing universal properties, we conclude that $F(A) \simeq F(\|A\|_0)$.

We can also use higher inductive types to construct colimits of algebraic objects. For instance, suppose $f : G \rightarrow H$ and $g : G \rightarrow K$ are group homomorphisms. Their pushout in the category of groups, called the **amalgamated free product** $H *_G K$, can be constructed as the higher inductive type generated by

- Functions $h : H \rightarrow H *_G K$ and $k : K \rightarrow H *_G K$.
- The operations and axioms of a group, as in the definition of $F(A)$.
- Axioms asserting that h and k are group homomorphisms.
- For $x : G$, we have $h(f(x)) = k(g(x))$.
- The 0-truncation constructor.

On the other hand, it can also be constructed explicitly, as the set-quotient of $\text{List}(H + K)$ by the following relations:

$$\begin{aligned} (\dots, x_1, x_2, \dots) &= (\dots, x_1 \cdot x_2, \dots) && \text{for } x_1, x_2 : H \\ (\dots, y_1, y_2, \dots) &= (\dots, y_1 \cdot y_2, \dots) && \text{for } y_1, y_2 : K \\ (\dots, 1_G, \dots) &= (\dots, \dots) \\ (\dots, 1_H, \dots) &= (\dots, \dots) \\ (\dots, f(x), \dots) &= (\dots, g(x), \dots) && \text{for } x : G. \end{aligned}$$

We leave the proofs to the reader. In the special case that G is the trivial group, the last relation is unnecessary, and we obtain the **free product** $H * K$, the coproduct in the category of groups. (This notation unfortunately clashes with that for the *join* of types, as in §6.8, but context generally disambiguates.)

Note that groups defined by *presentations* can be regarded as a special case of colimits. Suppose given a set (or more generally a type) A , and a pair of functions $R \Rightarrow F(A)$. We regard R as the type of “relations”, with the two functions assigning to each relation the two words that it sets equal. For instance, in the presentation $\langle a \mid a^2 = e \rangle$ we would have $A := \mathbf{1}$ and $R := \mathbf{1}$, with the two morphisms $R \Rightarrow F(A)$ picking out the list (a, a) and the empty list nil , respectively. Then by the universal property of free groups, we obtain a pair of group homomorphisms $F(R) \Rightarrow F(A)$. Their coequalizer in the category of groups, which can be built just like the pushout, is the group *presented* by this presentation.

Note that all these sorts of construction only apply to *algebraic* theories, which are theories whose axioms are (universally quantified) equations referring to variables, constants, and operations from a given signature. They can be modified to apply also to what are called *essentially algebraic theories*: those whose operations are partially defined on a domain specified by equalities between previous operations. They do not apply, for instance, to the theory of fields, in which the “inversion” operation is partially defined on a domain $\{x \mid x \neq 0\}$ specified by an *apartness* $\#$ between previous operations, see ???. And indeed, it is well-known that the category of fields has no initial object.

On the other hand, these constructions do apply just as well to *infinitary* algebraic theories, whose “operations” can take infinitely many inputs. In such cases, there may not be any presentation of free algebras or colimits of algebras as a simple quotient, unless we assume the axiom of choice. This means that higher inductive types represent a significant strengthening of constructive type theory (not necessarily in terms of proof-theoretic strength, but in terms of practical power), and indeed are stronger in some ways than Zermelo–Fraenkel set theory (without choice).

6.12 The flattening lemma

As we will see in ??, amazing things happen when we combine higher inductive types with univalence. The principal way this comes about is that if W is a higher inductive type and \mathcal{U} is a type universe, then we can define a type family $P : W \rightarrow \mathcal{U}$ by using the recursion principle for W . When we come to the clauses of the recursion principle dealing with the path constructors of W , we will need to supply paths in \mathcal{U} , and this is where univalence comes in.

For example, suppose we have a type X and a self-equivalence $e : X \simeq X$. Then we can

define a type family $P : \mathbb{S}^1 \rightarrow \mathcal{U}$ by using \mathbb{S}^1 -recursion:

$$P(\text{base}) := X \quad \text{and} \quad P(\text{loop}) := \text{ua}(e).$$

The type X thus appears as the fiber $P(\text{base})$ of P over the basepoint. The self-equivalence e is a little more hidden in P , but the following lemma says that it can be extracted by transporting along loop.

Lemma 6.12.1. *Given $B : A \rightarrow \mathcal{U}$ and $x, y : A$, with a path $p : x = y$ and an equivalence $e : B(x) \simeq B(y)$ such that $B(p) = \text{ua}(e)$, then for any $u : B(x)$ we have*

$$\text{transport}^B(p, u) = e(u).$$

Proof. Applying Lemma 2.10.5, we have

$$\begin{aligned} \text{transport}^B(p, u) &= \text{idtoeqv}(B(p))(u) \\ &= \text{idtoeqv}(\text{ua}(e))(u) \\ &= e(u). \end{aligned}$$

□

We have seen type families defined by recursion before: in §§2.12 and 2.13 we used them to characterize the identity types of (ordinary) inductive types. In ??, we will use similar ideas to calculate homotopy groups of higher inductive types.

In this section, we describe a general lemma about type families of this sort which will be useful later on. We call it the **flattening lemma**: it says that if $P : W \rightarrow \mathcal{U}$ is defined recursively as above, then its total space $\sum_{(x:W)} P(x)$ is equivalent to a “flattened” higher inductive type, whose constructors may be deduced from those of W and the definition of P . (From a category-theoretic point of view, $\sum_{(x:W)} P(x)$ is the “Grothendieck construction” of P , and the flattening lemma expresses its universal property as a “lax colimit”. Although because types in homotopy type theory (like W) correspond categorically to ∞ -groupoids (since all paths are invertible), in this case the lax colimit is the same as a pseudo colimit.)

We prove here one general case of the flattening lemma, which directly implies many particular cases and suggests the method to prove others. Suppose we have $A, B : \mathcal{U}$ and $f, g : B \rightarrow A$, and that the higher inductive type W is generated by

- $c : A \rightarrow W$ and
- $p : \prod_{(b:B)} (c(f(b)) =_W c(g(b)))$.

Thus, W is the **(homotopy) coequalizer** of f and g . Using binary sums (coproducts) and dependent sums (Σ -types), a lot of interesting nonrecursive higher inductive types can be represented in this form. All point constructors have to be bundled in the type A and all path constructors in the type B . For instance:

- The circle \mathbb{S}^1 can be represented by taking $A := \mathbf{1}$ and $B := \mathbf{1}$, with f and g the identity.
- The pushout of $j : X \rightarrow Y$ and $k : X \rightarrow Z$ can be represented by taking $A := Y + Z$ and $B := X$, with $f := \text{inl} \circ j$ and $g := \text{inr} \circ k$.

Now suppose in addition that

- $C : A \rightarrow \mathcal{U}$ is a family of types over A , and
- $D : \prod_{(b:B)} C(f(b)) \simeq C(g(b))$ is a family of equivalences over B .

Define a type family $P : W \rightarrow \mathcal{U}$ recursively by

$$\begin{aligned} P(c(a)) &:= C(a) \\ P(p(b)) &:= ua(D(b)). \end{aligned}$$

Let \tilde{W} be the higher inductive type generated by

- $\tilde{c} : \prod_{(a:A)} C(a) \rightarrow \tilde{W}$ and
- $\tilde{p} : \prod_{(b:B)} \prod_{(y:C(f(b)))} (\tilde{c}(f(b), y) =_{\tilde{W}} \tilde{c}(g(b), D(b)(y))).$

The flattening lemma is:

Lemma 6.12.2 (Flattening lemma). *In the above situation, we have*

$$\left(\sum_{x:W} P(x) \right) \simeq \tilde{W}.$$

As remarked above, this equivalence can be seen as expressing the universal property of $\sum_{(x:W)} P(x)$ as a “lax colimit” of P over W . It can also be seen as part of the *stability and descent* property of colimits, which characterizes higher toposes.

The proof of Lemma 6.12.2 occupies the rest of this section. It is somewhat technical and can be skipped on a first reading. But it is also a good example of “proof-relevant mathematics”, so we recommend it on a second reading.

The idea is to show that $\sum_{(x:W)} P(x)$ has the same universal property as \tilde{W} . We begin by showing that it comes with analogues of the constructors \tilde{c} and \tilde{p} .

Lemma 6.12.3. *There are functions*

- $\tilde{c}' : \prod_{(a:A)} C(a) \rightarrow \sum_{(x:W)} P(x)$ and
- $\tilde{p}' : \prod_{(b:B)} \prod_{(y:C(f(b)))} \left(\tilde{c}'(f(b), y) =_{\sum_{(w:W)} P(w)} \tilde{c}'(g(b), D(b)(y)) \right).$

Proof. The first is easy; define $\tilde{c}'(a, x) := (c(a), x)$ and note that by definition $P(c(a)) \equiv C(a)$. For the second, suppose given $b : B$ and $y : C(f(b))$; we must give an equality

$$(c(f(b)), y) = (c(g(b)), D(b)(y)).$$

Since we have $p(b) : c(f(b)) = c(g(b))$, by equalities in Σ -types it suffices to give an equality $p(b)_*(y) = D(b)(y)$. But this follows from Lemma 6.12.1, using the definition of P . \square

Now the following lemma says to define a section of a type family over $\sum_{(w:W)} P(w)$, it suffices to give analogous data as in the case of \tilde{W} .

Lemma 6.12.4. *Suppose $Q : (\sum_{(x:W)} P(x)) \rightarrow \mathcal{U}$ is a type family and that we have*

- $c : \prod_{(a:A)} \prod_{(x:C(a))} Q(\tilde{c}'(a, x))$ and
- $p : \prod_{(b:B)} \prod_{(y:C(f(b)))} \left(\tilde{p}'(b, y)_*(c(f(b), y)) = c(g(b), D(b)(y)) \right).$

Then there exists $k : \prod_{(z:\sum_{(w:W)} P(w))} Q(z)$ such that $k(\tilde{c}'(a, x)) \equiv c(a, x)$.

Proof. Suppose given $w : W$ and $x : P(w)$; we must produce an element $k(w, x) : Q(w, x)$. By induction on w , it suffices to consider two cases. When $w \equiv c(a)$, then we have $x : C(a)$, and so $c(a, x) : Q(c(a), x)$ as desired. (This part of the definition also ensures that the stated computational rule holds.)

Now we must show that this definition is preserved by transporting along $p(b)$ for any $b : B$. Since what we are defining, for all $w : W$, is a function of type $\prod_{(x:P(w))} Q(w, x)$, by Lemma 2.9.7 it suffices to show that for any $y : C(f(b))$, we have

$$\text{transport}^Q(\text{pair}^=(p(b), \text{refl}_{p(b)_*(y)}), c(f(b), y)) = c(g(b), p(b)_*(y)).$$

Let $q : p(b)_*(y) = D(b)(y)$ be the path obtained from Lemma 6.12.1. Then we have

$$\begin{aligned} c(g(b), p(b)_*(y)) &= \text{transport}^{x \mapsto Q(\tilde{c}'(g(b), x))}(q^{-1}, c(g(b), D(b)(y))) && (\text{by } \text{apd}_{x \mapsto c(g(b), x)}(q^{-1})^{-1}) \\ &= \text{transport}^Q(\text{ap}_{x \mapsto \tilde{c}'(g(b), x)}(q^{-1}), c(g(b), D(b)(y))). && (\text{by Lemma 2.3.10}) \end{aligned}$$

Thus, it suffices to show

$$\begin{aligned} \text{transport}^Q\left(\text{pair}^=(p(b), \text{refl}_{p(b)_*(y)}), c(f(b), y)\right) &= \\ \text{transport}^Q\left(\text{ap}_{x \mapsto \tilde{c}'(g(b), x)}(q^{-1}), c(g(b), D(b)(y))\right). \end{aligned}$$

Moving the right-hand transport to the other side, and combining two transports, this is equivalent to

$$\text{transport}^Q\left(\text{pair}^=(p(b), \text{refl}_{p(b)_*(y)}) \bullet \text{ap}_{x \mapsto \tilde{c}'(g(b), x)}(q), c(f(b), y)\right) = c(g(b), D(b)(y)).$$

However, we have

$$\begin{aligned} \text{pair}^=(p(b), \text{refl}_{p(b)_*(y)}) \bullet \text{ap}_{x \mapsto \tilde{c}'(g(b), x)}(q) &= \\ \text{pair}^=(p(b), \text{refl}_{p(b)_*(y)}) \bullet \text{pair}^=(\text{refl}_{c(g(b))}, q) &= \text{pair}^=(p(b), q) = \tilde{p}'(b, y) \end{aligned}$$

so the construction is completed by the assumption $p(b, y)$ of type

$$\text{transport}^Q(\tilde{p}'(b, y), c(f(b), y)) = c(g(b), D(b)(y)). \quad \square$$

Lemma 6.12.4 almost gives $\sum_{(w:W)} P(w)$ the same induction principle as \tilde{W} . The missing bit is the equality $\text{apd}_k(\tilde{p}'(b, y)) = p(b, y)$. In order to prove this, we would need to analyze the proof of Lemma 6.12.4, which of course is the definition of k .

It should be possible to do this, but it turns out that we only need the computation rule for the non-dependent recursion principle. Thus, we now give a somewhat simpler direct construction of the recursor, and a proof of its computation rule.

Lemma 6.12.5. *Suppose Q is a type and that we have*

- $c : \prod_{(a:A)} C(a) \rightarrow Q$ and
- $p : \prod_{(b:B)} \prod_{(y:C(f(b)))} (c(f(b), y) =_Q c(g(b), D(b)(y)))$.

Then there exists $k : (\sum_{(w:W)} P(w)) \rightarrow Q$ such that $k(\tilde{c}'(a, x)) \equiv c(a, x)$.

Proof. As in Lemma 6.12.4, we define $k(w, x)$ by induction on $w : W$. When $w \equiv c(a)$, we define $k(c(a), x) := c(a, x)$. Now by Lemma 2.9.6, it suffices to consider, for $b : B$ and $y : C(f(b))$, the composite path

$$\text{transport}^{x \mapsto Q}(\mathbf{p}(b), c(f(b), y)) = c(g(b), \text{transport}^P(\mathbf{p}(b), y)) \quad (6.12.6)$$

defined as the composition

$$\begin{aligned} \text{transport}^{x \mapsto Q}(\mathbf{p}(b), c(f(b), y)) &= c(f(b), y) && \text{(by Lemma 2.3.5)} \\ &= c(g(b), D(b)(y)) && \text{(by } p(b, y)\text{)} \\ &= c(g(b), \text{transport}^P(\mathbf{p}(b), y)). && \text{(by Lemma 6.12.1)} \end{aligned}$$

The computation rule $k(\tilde{c}'(a, x)) \equiv c(a, x)$ follows by definition, as before. \square

For the second computation rule, we need the following lemma.

Lemma 6.12.7. *Let $Y : X \rightarrow \mathcal{U}$ be a type family and let $k : (\sum_{(x:X)} Y(x)) \rightarrow Z$ be defined componentwise by $k(x, y) := d(x)(y)$ for a curried function $d : \prod_{(x:X)} Y(x) \rightarrow Z$. Then for any $s : x_1 =_X x_2$ and any $y_1 : Y(x_1)$ and $y_2 : Y(x_2)$ with a path $r : s_*(y_1) = y_2$, the path*

$$\text{ap}_k(\text{pair}^=(s, r)) : k(x_1, y_1) = k(x_2, y_2)$$

is equal to the composite

$$\begin{aligned} k(x_1, y_1) &\equiv d(x_1)(y_1) \\ &= \text{transport}^{x \mapsto Z}(s, d(x_1)(y_1)) && \text{(by (Lemma 2.3.5)}^{-1}\text{)} \\ &= \text{transport}^{x \mapsto Z}(s, d(x_1)(s^{-1}_*(s_*(y_1)))) \\ &= (\text{transport}^{x \mapsto (Y(x) \rightarrow Z)}(s, d(x_1))(s_*(y_1))) && \text{(by (2.9.4))} \\ &= d(x_2)(s_*(y_1)) && \text{(by } \text{happly}(\text{apd}_d(s))(s_*(y_1))\text{)} \\ &= d(x_2)(y_2) && \text{(by } \text{ap}_{d(x_2)}(r)\text{)} \\ &\equiv k(x_2, y_2). \end{aligned}$$

Proof. After path induction on s and r , both equalities reduce to reflexivities. \square

At first it may seem surprising that Lemma 6.12.7 has such a complicated statement, while it can be proven so simply. The reason for the complication is to ensure that the statement is well-typed: $\text{ap}_k(\text{pair}^=(s, r))$ and the composite path it is claimed to be equal to must both have the same start and end points. Once we have managed this, the proof is easy by path induction.

Lemma 6.12.8. *In the situation of Lemma 6.12.5, we have $\text{ap}_k(\tilde{p}'(b, y)) = p(b, y)$.*

Proof. Recall that $\tilde{p}'(b, y) := \text{pair}^=(\mathbf{p}(b), q)$ where $q : \mathbf{p}(b)_*(y) = D(b)(y)$ comes from Lemma 6.12.1. Thus, since k is defined componentwise, we may compute $\text{ap}_k(\tilde{p}'(b, y))$ by Lemma 6.12.7, with

$$\begin{array}{ll} x_1 := c(f(b)) & y_1 := y \\ x_2 := c(g(b)) & y_2 := D(b)(y) \\ s := \mathbf{p}(b) & r := q. \end{array}$$

The curried function $d : \prod_{(w:W)} P(w) \rightarrow Q$ was defined by induction on $w : W$; to apply Lemma 6.12.7 we need to understand $\text{ap}_{d(x_2)}(r)$ and $\text{happly}(\text{apd}_d(s), s_*(y_1))$.

For the first, since $d(c(a), x) \equiv c(a, x)$, we have

$$\text{ap}_{d(x_2)}(r) \equiv \text{ap}_{c(g(b), -)}(q).$$

For the second, the computation rule for the induction principle of W tells us that $\text{apd}_d(p(b))$ is equal to the composite (6.12.6), passed across the equivalence of Lemma 2.9.6. Thus, the computation rule given in Lemma 2.9.6 implies that $\text{happly}(\text{apd}_d(p(b)), p(b)_*(y))$ is equal to the composite

$$\begin{aligned} (p(b)_*(c(f(b), -)))(p(b)_*(y)) &= p(b)_*(c(f(b), p(b)^{-1}_*(p(b)_*(y)))) && (\text{by (2.9.4)}) \\ &= p(b)_*(c(f(b), y)) \\ &= c(f(b), y) && (\text{by Lemma 2.3.5}) \\ &= c(f(b), D(b)(y)) && (\text{by } p(b, y)) \\ &= c(f(b), p(b)_*(y)). && (\text{by } \text{ap}_{c(g(b), -)}(q)^{-1}) \end{aligned}$$

Finally, substituting these values of $\text{ap}_{d(x_2)}(r)$ and $\text{happly}(\text{apd}_d(s), s_*(y_1))$ into Lemma 6.12.7, we see that all the paths cancel out in pairs, leaving only $p(b, y)$. \square

Now we are finally ready to prove the flattening lemma.

Proof of Lemma 6.12.2. We define $h : \tilde{W} \rightarrow \sum_{(w:W)} P(w)$ by using the recursion principle for \tilde{W} , with \tilde{c}' and \tilde{p}' as input data. Similarly, we define $k : (\sum_{(w:W)} P(w)) \rightarrow \tilde{W}$ by using the recursion principle of Lemma 6.12.5, with \tilde{c} and \tilde{p} as input data.

On the one hand, we must show that for any $z : \tilde{W}$, we have $k(h(z)) = z$. By induction on z , it suffices to consider the two constructors of \tilde{W} . But we have

$$k(h(\tilde{c}(a, x))) \equiv k(\tilde{c}'(a, x)) \equiv \tilde{c}(a, x)$$

by definition, while similarly

$$k(h(\tilde{p}(b, y))) = k(\tilde{p}'(b, y)) = \tilde{p}(b, y)$$

using the propositional computation rule for \tilde{W} and Lemma 6.12.8.

On the other hand, we must show that for any $z : \sum_{(w:W)} P(w)$, we have $h(k(z)) = z$. But this is essentially identical, using Lemma 6.12.4 for “induction on $\sum_{(w:W)} P(w)$ ” and the same computation rules. \square

6.13 The general syntax of higher inductive definitions

In §5.6, we discussed the conditions on a putative “inductive definition” which make it acceptable, namely that all inductive occurrences of the type in its constructors are “strictly positive”. In this section, we say something about the additional conditions required for *higher* inductive definitions. Finding a general syntactic description of valid higher inductive definitions is an area of current research, and all of the solutions proposed to date are somewhat technical in nature; thus we only give a general description and not a precise definition. Fortunately, the corner cases never seem to arise in practice.

Like an ordinary inductive definition, a higher inductive definition is specified by a list of *constructors*, each of which is a (dependent) function. For simplicity, we may require the inputs

of each constructor to satisfy the same condition as the inputs for constructors of ordinary inductive types. In particular, they may contain the type being defined only strictly positively. Note that this excludes definitions such as the 0-truncation as presented in §6.9, where the input of a constructor contains not only the inductive type being defined, but its identity type as well. It may be possible to extend the syntax to allow such definitions; but also, in §7.3 we will give a different construction of the 0-truncation whose constructors do satisfy the more restrictive condition.

The only difference between an ordinary inductive definition and a higher one, then, is that the *output* type of a constructor may be, not the type being defined (W , say), but some identity type of it, such as $u =_W v$, or more generally an iterated identity type such as $p =_{(u=_W v)} q$. Thus, when we give a higher inductive definition, we have to specify not only the inputs of each constructor, but the expressions u and v (or u, v, p , and q , etc.) which determine the source and target of the path being constructed.

Importantly, these expressions may refer to *other* constructors of W . For instance, in the definition of S^1 , the constructor `loop` has both u and v being base, the previous constructor. To make sense of this, we require the constructors of a higher inductive type to be specified *in order*, and we allow the source and target expressions u and v of each constructor to refer to previous constructors, but not later ones. (Of course, in practice the constructors of any inductive definition are written down in some order, but for ordinary inductive types that order is irrelevant.)

Note that this order is not necessarily the order of “dimension”: in principle, a 1-dimensional path constructor could refer to a 2-dimensional one and hence need to come after it. However, we have not given the 0-dimensional constructors (point constructors) any way to refer to previous constructors, so they might as well all come first. And if we use the hub-and-spoke construction (§6.7) to reduce all constructors to points and 1-paths, then we might assume that all point constructors come first, followed by all 1-path constructors — but the order among the 1-path constructors continues to matter.

The remaining question is, what sort of expressions can u and v be? We might hope that they could be any expression at all involving the previous constructors. However, the following example shows that a naive approach to this idea does not work.

Example 6.13.1. Consider a family of functions $f : \prod_{(X:\mathcal{U})}(X \rightarrow X)$. Of course, f_X might be just id_X for all X , but other such f s may also exist. For instance, nothing prevents $f_2 : 2 \rightarrow 2$ from being the nonidentity automorphism (see Exercise 6.9).

Now suppose that we attempt to define a higher inductive type K generated by:

- two elements $a, b : K$, and
- a path $\sigma : f_K(a) = f_K(b)$.

What would the induction principle for K say? We would assume a type family $P : K \rightarrow \mathcal{U}$, and of course we would need $x : P(a)$ and $y : P(b)$. The remaining datum should be a dependent path in P living over σ , which must therefore connect some element of $P(f_K(a))$ to some element of $P(f_K(b))$. But what could these elements possibly be? We know that $P(a)$ and $P(b)$ are inhabited by x and y , respectively, but this tells us nothing about $P(f_K(a))$ and $P(f_K(b))$.

Clearly some condition on u and v is required in order for the definition to be sensible. It seems that, just as the domain of each constructor is required to be (among other things) a *covariant functor*, the appropriate condition on the expressions u and v is that they define *natural transformations*. Making precise sense of this requirement is beyond the scope of this book, but informally it means that u and v must only involve operations which are preserved by all functions between types.

For instance, it is permissible for u and v to refer to concatenation of paths, as in the case of the final constructor of the torus in §6.6, since all functions in type theory preserve path concatenation (up to homotopy). However, it is not permissible for them to refer to an operation like the function f in Example 6.13.1, which is not necessarily natural: there might be some function $g : X \rightarrow Y$ such that $f_Y \circ g \neq g \circ f_X$. (Univalence implies that f_X must be natural with respect to all *equivalences*, but not necessarily with respect to functions that are not equivalences.)

The intuition of naturality supplies only a rough guide for when a higher inductive definition is permissible. Even if it were possible to give a precise specification of permissible forms of such definitions in this book, such a specification would probably be out of date quickly, as new extensions to the theory are constantly being explored. For instance, the presentation of n -spheres in terms of “dependent n -loops” referred to in §6.4, and the “higher inductive-recursive definitions” used in ??, were innovations introduced while this book was being written. We encourage the reader to experiment — with caution.

Notes

The general idea of higher inductive types was conceived in discussions between Andrej Bauer, Peter Lumsdaine, Mike Shulman, and Michael Warren at the Oberwolfach meeting in 2011, although there are some suggestions of some special cases in earlier work. Subsequently, Guillaume Brunerie and Dan Licata contributed substantially to the general theory, especially by finding convenient ways to represent them in computer proof assistants and do homotopy theory with them (see ??).

A general discussion of the syntax of higher inductive types, and their semantics in higher-categorical models, appears in [LS17]. As with ordinary inductive types, models of higher inductive types can be constructed by transfinite iterative processes; a slogan is that ordinary inductive types describe *free* monads while higher inductive types describe *presentations* of monads. The introduction of path constructors also involves the model-category-theoretic equivalence between “right homotopies” (defined using path spaces) and “left homotopies” (defined using cylinders) — the fact that this equivalence is generally only up to homotopy provides a semantic reason to prefer propositional computation rules for path constructors.

Another (temporary) reason for this preference comes from the limitations of existing computer implementations. Proof assistants like COQ and AGDA have ordinary inductive types built in, but not yet higher inductive types. We can of course introduce them by assuming lots of axioms, but this results in only propositional computation rules. However, there is a trick due to Dan Licata which implements higher inductive types using private data types; this yields judgmental rules for point constructors but not path constructors.

The type-theoretic description of higher spheres using loop spaces and suspensions in §§6.4 and 6.5 is largely due to Brunerie and Licata; Hou has given a type-theoretic version of the alternative description that uses n -dimensional paths. The reduction of higher paths to 1-dimensional paths with hubs and spokes (§6.7) is due to Lumsdaine and Shulman. The description of truncation as a higher inductive type is due to Lumsdaine; the (-1) -truncation is closely related to the “bracket types” of [AB04]. The flattening lemma was first formulated in generality by Brunerie.

Quotient types are unproblematic in extensional type theory, such as NUPRL [CAB⁺86]. They are often added by passing to an extended system of setoids. However, quotients are a trickier issue in intensional type theory (the starting point for homotopy type theory), because one cannot simply add new propositional equalities without specifying how they are to behave.

Some solutions to this problem have been studied [Hof95, Alt99, AMS07], and several different notions of quotient types have been considered. The construction of set-quotients using higher-inductives provides an argument for our particular approach (which is similar to some that have previously been considered), because it arises as an instance of a general mechanism. Our construction does not yet provide a new solution to all the computational problems related to quotients, since we still lack a good computational understanding of higher inductive types in general—but it does mean that ongoing work on the computational interpretation of higher inductives applies to the quotients as well. The construction of quotients in terms of equivalence classes is, of course, a standard set-theoretic idea, and a well-known aspect of elementary topos theory; its use in type theory (which depends on the univalence axiom, at least for mere propositions) was proposed by Voevodsky. The fact that quotient types in intensional type theory imply function extensionality was proved by [Hof95], inspired by the work of [Car95] on exact completions; Lemma 6.3.2 is an adaptation of such arguments.

Exercises

Exercise 6.1. Define concatenation of dependent paths, prove that application of dependent functions preserves concatenation, and write out the precise induction principle for the torus T^2 with its computation rules.

Exercise 6.2. Prove that $\Sigma \mathbb{S}^1 \simeq \mathbb{S}^2$, using the explicit definition of \mathbb{S}^2 in terms of base and surf given in §6.4.

Exercise 6.3. Prove that the torus T^2 as defined in §6.6 is equivalent to $\mathbb{S}^1 \times \mathbb{S}^1$. (Warning: the path algebra for this is rather difficult.)

Exercise 6.4. Define dependent n -loops and the action of dependent functions on n -loops, and write down the induction principle for the n -spheres as defined at the end of §6.4.

Exercise 6.5. Prove that $\Sigma \mathbb{S}^n \simeq \mathbb{S}^{n+1}$, using the definition of \mathbb{S}^n in terms of Ω^n from §6.4.

Exercise 6.6. Prove that if the type \mathbb{S}^2 belongs to some universe \mathcal{U} , then \mathcal{U} is not a 2-type.

Exercise 6.7. Prove that if G is a monoid and $x : G$, then $\sum_{(y:G)} ((x \cdot y = e) \times (y \cdot x = e))$ is a mere proposition. Conclude, using the principle of unique choice (Corollary 3.9.2), that it would be equivalent to define a group to be a monoid such that for every $x : G$, there merely exists a $y : G$ such that $x \cdot y = e$ and $y \cdot x = e$.

Exercise 6.8. Prove that if A is a set, then $\text{List}(A)$ is a monoid. Then complete the proof of Lemma 6.11.5.

Exercise 6.9. Assuming LEM, construct a family $f : \prod_{(X:\mathcal{U})} (X \rightarrow X)$ such that $f_2 : \mathbf{2} \rightarrow \mathbf{2}$ is the nonidentity automorphism.

Exercise 6.10. Show that the map constructed in Lemma 6.3.2 is in fact a quasi-inverse to happily, so that an interval type implies the full function extensionality axiom. (You may have to use Exercise 2.16.)

Exercise 6.11. Prove the universal property of suspension:

$$(\Sigma A \rightarrow B) \simeq \left(\sum_{(b_n:B)} \sum_{(b_s:B)} (A \rightarrow (b_n = b_s)) \right)$$

Exercise 6.12. Show that $\mathbb{Z} \simeq \mathbb{N} + \mathbf{1} + \mathbb{N}$. Show that if we were to define \mathbb{Z} as $\mathbb{N} + \mathbf{1} + \mathbb{N}$, then we could obtain Lemma 6.10.12 with judgmental computation rules.

Exercise 6.13. Show that we can also prove Lemma 6.3.2 by using $\|\mathbf{2}\|$ instead of I .

Chapter 7

Homotopy n -types

One of the basic notions of homotopy theory is that of a *homotopy n -type*: a space containing no interesting homotopy above dimension n . For instance, a homotopy 0-type is essentially a set, containing no nontrivial paths, while a homotopy 1-type may contain nontrivial paths, but no nontrivial paths between paths. Homotopy n -types are also called *n -truncated spaces*. We have mentioned this notion already in §3.1; our first goal in this chapter is to give it a precise definition in homotopy type theory.

A dual notion to truncatedness is connectedness: a space is *n -connected* if it has no interesting homotopy in dimensions n and *below*. For instance, a space is 0-connected (also called just “connected”) if it has only one connected component, and 1-connected (also called “simply connected”) if it also has no nontrivial loops (though it may have nontrivial higher loops between loops).

The duality between truncatedness and connectedness is most easily seen by extending both notions to maps. We call a map *n -truncated* or *n -connected* if all its fibers are so. Then n -connected and n -truncated maps form the two classes of maps in an *orthogonal factorization system*, i.e. every map factors uniquely as an n -connected map followed by an n -truncated one.

In the case $n = -1$, the n -truncated maps are the embeddings and the n -connected maps are the surjections, as defined in §4.6. Thus, the n -connected factorization system is a massive generalization of the standard image factorization of a function between sets into a surjection followed by an injection. At the end of this chapter, we sketch briefly an even more general theory: any type-theoretic *modality* gives rise to an analogous factorization system.

7.1 Definition of n -types

As mentioned in §§3.1 and 3.11, it turns out to be convenient to define n -types starting two levels below zero, with the (-1) -types being the mere propositions and the (-2) -types the contractible ones.

Definition 7.1.1. Define the predicate $\text{is-}n\text{-type} : \mathcal{U} \rightarrow \mathcal{U}$ for $n \geq -2$ by recursion as follows:

$$\text{is-}n\text{-type}(X) := \begin{cases} \text{isContr}(X) & \text{if } n = -2, \\ \prod_{(x,y:X)} \text{is-}n'\text{-type}(x =_X y) & \text{if } n = n' + 1. \end{cases}$$

We say that X is an *n -type*, or sometimes that it is *n -truncated*, if $\text{is-}n\text{-type}(X)$ is inhabited.

Remark 7.1.2. The number n in Definition 7.1.1 ranges over all integers greater than or equal to -2 . We could make sense of this formally by defining a type $\mathbb{Z}_{\geq -2}$ of such integers (a type whose induction principle is identical to that of \mathbb{N}), or instead defining a predicate $\text{is-}(k-2)$ -type for $k : \mathbb{N}$. Either way, we can prove theorems about n -types by induction on n , with $n = -2$ as the base case.

Example 7.1.3. We saw in Lemma 3.11.10 that X is a (-1) -type if and only if it is a mere proposition. Therefore, X is a 0 -type if and only if it is a set.

We have also seen that there are types which are not sets (Example 3.1.9). So far, however, we have not shown for any $n > 0$ that there exist types which are not n -types. In ??, however, we will show that the $(n+1)$ -sphere S^{n+1} is not an n -type. (Kraus has also shown that the n^{th} nested univalent universe is also not an n -type, without using any higher inductive types.) Moreover, in ?? will give an example of a type that is not an n -type for *any* (finite) number n .

We begin the general theory of n -types by showing they are closed under certain operations and constructors.

Theorem 7.1.4. *Let $p : X \rightarrow Y$ be a retraction and suppose that X is an n -type, for any $n \geq -2$. Then Y is also an n -type.*

Proof. We proceed by induction on n . The base case $n = -2$ is handled by Lemma 3.11.7.

For the inductive step, assume that any retract of an n -type is an n -type, and that X is an $(n+1)$ -type. Let $y, y' : Y$; we must show that $y = y'$ is an n -type. Let s be a section of p , and let ϵ be a homotopy $\epsilon : p \circ s \sim 1$. Since X is an $(n+1)$ -type, $s(y) =_X s(y')$ is an n -type. We claim that $y = y'$ is a retract of $s(y) =_X s(y')$. For the section, we take

$$\text{ap}_s : (y = y') \rightarrow (s(y) = s(y')).$$

For the retraction, we define $t : (s(y) = s(y')) \rightarrow (y = y')$ by

$$t(q) := \epsilon_y^{-1} \cdot p(q) \cdot \epsilon_{y'}.$$

To show that t is a retraction of ap_s , we must show that

$$\epsilon_y^{-1} \cdot p(s(r)) \cdot \epsilon_{y'} = r$$

for any $r : y = y'$. But this follows from Lemma 2.4.3. \square

As an immediate corollary we obtain the stability of n -types under equivalence (which is also immediate from univalence):

Corollary 7.1.5. *If $X \simeq Y$ and X is an n -type, then so is Y .*

Recall also the notion of embedding from §4.6.

Theorem 7.1.6. *If $f : X \rightarrow Y$ is an embedding and Y is an n -type for some $n \geq -1$, then so is X .*

Proof. Let $x, x' : X$; we must show that $x =_X x'$ is an $(n-1)$ -type. But since f is an embedding, we have $(x =_X x') \simeq (f(x) =_Y f(x'))$, and the latter is an $(n-1)$ -type by assumption. \square

Note that this theorem fails when $n = -2$: the map $\mathbf{0} \rightarrow \mathbf{1}$ is an embedding, but $\mathbf{1}$ is a (-2) -type while $\mathbf{0}$ is not.

Theorem 7.1.7. *The hierarchy of n -types is cumulative in the following sense: given a number $n \geq -2$, if X is an n -type, then it is also an $(n + 1)$ -type.*

Proof. We proceed by induction on n .

For $n = -2$, we need to show that a contractible type, say, A , has contractible path spaces. Let $a_0 : A$ be the center of contraction of A , and let $x, y : A$. We show that $x =_A y$ is contractible. By contractibility of A we have a path $\text{contr}_x \bullet \text{contr}_y^{-1} : x = y$, which we choose as the center of contraction for $x = y$. Given any $p : x = y$, we need to show $p = \text{contr}_x \bullet \text{contr}_y^{-1}$. By path induction, it suffices to show that $\text{refl}_x = \text{contr}_x \bullet \text{contr}_x^{-1}$, which is trivial.

For the inductive step, we need to show that $x =_X y$ is an $(n + 1)$ -type, provided that X is an $(n + 1)$ -type. Applying the inductive hypothesis to $x =_X y$ yields the desired result. \square

We now show that n -types are preserved by most of the type forming operations.

Theorem 7.1.8. *Let $n \geq -2$, and let $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$. If A is an n -type and for all $a : A$, $B(a)$ is an n -type, then so is $\sum_{(x:A)} B(x)$.*

Proof. We proceed by induction on n .

For $n = -2$, we choose the center of contraction for $\sum_{(x:A)} B(x)$ to be the pair (a_0, b_0) , where $a_0 : A$ is the center of contraction of A and $b_0 : B(a_0)$ is the center of contraction of $B(a_0)$. Given any other element (a, b) of $\sum_{(x:A)} B(x)$, we provide a path $(a, b) = (a_0, b_0)$ by contractibility of A and $B(a_0)$, respectively.

For the inductive step, suppose that A is an $(n + 1)$ -type and for any $a : A$, $B(a)$ is an $(n + 1)$ -type. We show that $\sum_{(x:A)} B(x)$ is an $(n + 1)$ -type: fix (a_1, b_1) and (a_2, b_2) in $\sum_{(x:A)} B(x)$, we show that $(a_1, b_1) = (a_2, b_2)$ is an n -type. By Theorem 2.7.2 we have

$$((a_1, b_1) = (a_2, b_2)) \simeq \sum_{p:a_1=a_2} (p_*(b_1) =_{B(a_2)} b_2)$$

and by preservation of n -types under equivalences (Corollary 7.1.5) it suffices to prove that the latter is an n -type. This follows from the inductive hypothesis. \square

As a special case, if A and B are n -types, so is $A \times B$. Note also that Theorem 7.1.7 implies that if A is an n -type, then so is $x =_A y$ for any $x, y : A$. Combining this with Theorem 7.1.8, we see that for any functions $f : A \rightarrow C$ and $g : B \rightarrow C$ between n -types, their pullback

$$A \times_C B := \sum_{(x:A)} \sum_{(y:B)} (f(x) = g(y))$$

(see Exercise 2.11) is also an n -type. More generally, n -types are closed under all *limits*.

Theorem 7.1.9. *Let $n \geq -2$, and let $A : \mathcal{U}$ and $B : A \rightarrow \mathcal{U}$. If for all $a : A$, $B(a)$ is an n -type, then so is $\prod_{(x:A)} B(x)$.*

Proof. We proceed by induction on n . For $n = -2$, the result is simply Lemma 3.11.6.

For the inductive step, assume the result is true for n -types, and that each $B(a)$ is an $(n + 1)$ -type. Let $f, g : \prod_{(a:A)} B(a)$. We need to show that $f = g$ is an n -type. By function extensionality and closure of n -types under equivalence, it suffices to show that $\prod_{(a:A)} (f(a) =_{B(a)} g(a))$ is an n -type. This follows from the inductive hypothesis. \square

As a special case of the above theorem, the function space $A \rightarrow B$ is an n -type provided that B is an n -type. We can now generalize our observations in Chapter 2 that $\text{isSet}(A)$ and $\text{isProp}(A)$ are mere propositions.

Theorem 7.1.10. *For any $n \geq -2$ and any type X , the type $\text{is-}n\text{-type}(X)$ is a mere proposition.*

Proof. We proceed by induction with respect to n .

For the base case, we need to show that for any X , the type $\text{isContr}(X)$ is a mere proposition. This is Lemma 3.11.4.

For the inductive step we need to show

$$\prod_{X:\mathcal{U}} \text{isProp}(\text{is-}n\text{-type}(X)) \rightarrow \prod_{X:\mathcal{U}} \text{isProp}(\text{is-}(n+1)\text{-type}(X)).$$

To show the conclusion of this implication, we need to show that for any type X , the type

$$\prod_{x,x':X} \text{is-}n\text{-type}(x = x')$$

is a mere proposition. By Example 3.6.2 or Theorem 7.1.9, it suffices to show that for any $x, x' : X$, the type $\text{is-}n\text{-type}(x =_X x')$ is a mere proposition. But this follows from the inductive hypothesis applied to the type $(x =_X x')$. \square

Finally, we show that the type of n -types is itself an $(n+1)$ -type. We define this to be:

$$n\text{-Type} := \sum_{X:\mathcal{U}} \text{is-}n\text{-type}(X).$$

If necessary, we may specify the universe \mathcal{U} by writing $n\text{-Type}_{\mathcal{U}}$. In particular, we have $\text{Prop} := (-1)\text{-Type}$ and $\text{Set} := 0\text{-Type}$, as defined in Chapter 2. Note that just as for Prop and Set , because $\text{is-}n\text{-type}(X)$ is a mere proposition, by Lemma 3.5.1 for any $(X, p), (X', p') : n\text{-Type}$ we have

$$\begin{aligned} ((X, p) =_{n\text{-Type}} (X', p')) &\simeq (X =_{\mathcal{U}} X') \\ &\simeq (X \simeq X'). \end{aligned}$$

Theorem 7.1.11. *For any $n \geq -2$, the type $n\text{-Type}$ is an $(n+1)$ -type.*

Proof. Let $(X, p), (X', p') : n\text{-Type}$; we need to show that $(X, p) = (X', p')$ is an n -type. By the above observation, this type is equivalent to $X \simeq X'$. Next, we observe that the projection

$$(X \simeq X') \rightarrow (X \rightarrow X').$$

is an embedding, so that if $n \geq -1$, then by Theorem 7.1.6 it suffices to show that $X \rightarrow X'$ is an n -type. But since n -types are preserved under the arrow type, this reduces to an assumption that X' is an n -type.

In the case $n = -2$, this argument shows that $X \simeq X'$ is a (-1) -type—but it is also inhabited, since any two contractible types are equivalent to $\mathbf{1}$, and hence to each other. Thus, $X \simeq X'$ is also a (-2) -type. \square

7.2 Uniqueness of identity proofs and Hedberg's theorem

In §3.1 we defined a type X to be a *set* if for all $x, y : X$ and $p, q : x =_X y$ we have $p = q$. In conventional type theory, this property goes by the name of **uniqueness of identity proofs (UIP)**. We have seen also that it is equivalent to being a 0-type in the sense of the previous section. Here is another equivalent characterization, involving Streicher's “Axiom K” [Str93]:

Theorem 7.2.1. *A type X is a set if and only if it satisfies Axiom K: for all $x : X$ and $p : (x =_A x)$ we have $p = \text{refl}_x$.*

Proof. Clearly Axiom K is a special case of UIP. Conversely, if X satisfies Axiom K, let $x, y : X$ and $p, q : (x = y)$; we want to show $p = q$. But induction on q reduces this goal precisely to Axiom K. \square

We stress that we are not assuming UIP or the K principle as axioms! They are simply properties which a particular type may or may not satisfy (which are equivalent to being a set). Recall from Example 3.1.9 that *not* all types are sets.

The following theorem is another useful way to show that types are sets.

Theorem 7.2.2. *Suppose R is a reflexive mere relation on a type X implying identity. Then X is a set, and $R(x, y)$ is equivalent to $x =_X y$ for all $x, y : X$.*

Proof. Let $\rho : \prod_{(x:X)} R(x, x)$ witness reflexivity of R , and let $f : \prod_{(x,y:X)} R(x, y) \rightarrow (x =_X y)$ be a witness that R implies identity. Note first that the two statements in the theorem are equivalent. For on one hand, if X is a set, then $x =_X y$ is a mere proposition, and since it is logically equivalent to the mere proposition $R(x, y)$ by hypothesis, it must also be equivalent to it. On the other hand, if $x =_X y$ is equivalent to $R(x, y)$, then like the latter it is a mere proposition for all $x, y : X$, and hence X is a set.

We give two proofs of this theorem. The first shows directly that X is a set; the second shows directly that $R(x, y) \simeq (x = y)$.

First proof: we show that X is a set. The idea is the same as that of Lemma 3.3.4: the function f must be continuous in its arguments x and y . However, it is slightly more notationally complicated because we have to deal with the additional argument of type $R(x, y)$.

Firstly, for any $x : X$ and $p : x =_X x$, consider $\text{apd}_{f(x)}(p)$. This is a dependent path from $f(x, x)$ to itself. Since $f(x, x)$ is still a function $R(x, x) \rightarrow (x =_X x)$, by Lemma 2.9.6 this yields for any $r : R(x, x)$ a path

$$p_*(f(x, x, r)) = f(x, x, p_*(r)).$$

On the left-hand side, we have transport in an identity type, which is concatenation. And on the right-hand side, we have $p_*(r) = r$, since both lie in the mere proposition $R(x, x)$. Thus, substituting $r := \rho(x)$, we obtain

$$f(x, x, \rho(x)) \cdot p = f(x, x, \rho(x)).$$

By cancellation, $p = \text{refl}_x$. So X satisfies Axiom K, and hence is a set.

Second proof: we show that each $f(x, y) : R(x, y) \rightarrow x =_X y$ is an equivalence. By Theorem 4.7.7, it suffices to show that f induces an equivalence of total spaces:

$$\left(\sum_{y:X} R(x, y) \right) \simeq \left(\sum_{y:X} x =_X y \right).$$

By Lemma 3.11.8, the type on the right is contractible, so it suffices to show that the type on the left is contractible. As the center of contraction we take the pair $(x, \rho(x))$. It remains to show, for every $y : X$ and every $H : R(x, y)$ that

$$(x, \rho(x)) = (y, H).$$

But since $R(x, y)$ is a mere proposition, by Theorem 2.7.2 it suffices to show that $x =_X y$, which we get from $f(H)$. \square

Corollary 7.2.3. *If a type X has the property that $\neg\neg(x = y) \rightarrow (x = y)$ for any $x, y : X$, then X is a set.*

Another convenient way to show that a type is a set is the following. Recall from §3.4 that a type X is said to have *decidable equality* if for all $x, y : X$ we have

$$(x =_X y) + \neg(x =_X y).$$

This is a very strong condition: it says that a path $x = y$ can be chosen, when it exists, continuously (or computably, or functorially) in x and y . This turns out to imply that X is a set, by way of Theorem 7.2.2 and the following lemma.

Lemma 7.2.4. *For any type A we have $(A + \neg A) \rightarrow (\neg\neg A \rightarrow A)$.*

Proof. This was essentially already proven in Corollary 3.2.7, but we repeat the argument. Suppose $x : A + \neg A$. We have two cases to consider. If x is $\text{inl}(a)$ for some $a : A$, then we have the constant function $\neg\neg A \rightarrow A$ which maps everything to a . If x is $\text{inr}(t)$ for some $t : \neg A$, we have $g(t) : \mathbf{0}$ for every $g : \neg\neg A$. Hence we may use *ex falso quodlibet*, that is rec_0 , to obtain an element of A for any $g : \neg\neg A$. \square

Theorem 7.2.5 (Hedberg). *If X has decidable equality, then X is a set.*

Proof. If X has decidable equality, it follows that $\neg\neg(x = y) \rightarrow (x = y)$ for any $x, y : X$. Therefore, Hedberg's theorem follows from Corollary 7.2.3. \square

There is, of course, a strong connection between this theorem and Corollary 3.2.7. The statement LEM_∞ that is denied by Corollary 3.2.7 clearly implies that every type has decidable equality, and hence is a set, which we know is not the case. Note that the consistent axiom LEM from §3.4 implies only that every type has *merely decidable equality*, i.e. that for any A we have

$$\prod_{a,b:A} (\|a = b\| + \neg\|a = b\|).$$

As an example application of Theorem 7.2.5, recall that in Example 3.1.4 we observed that \mathbb{N} is a set, using our characterization of its equality types in §2.13. A more traditional proof of this theorem uses only (2.13.2) and (2.13.3), rather than the full characterization of Theorem 2.13.1, with Theorem 7.2.5 to fill in the blanks.

Theorem 7.2.6. *The type \mathbb{N} of natural numbers has decidable equality, and hence is a set.*

Proof. Let $x, y : \mathbb{N}$ be given; we proceed by induction on x and case analysis on y to prove $(x = y) + \neg(x = y)$. If $x \equiv 0$ and $y \equiv 0$, we take $\text{inl}(\text{refl}_0)$. If $x \equiv 0$ and $y \equiv \text{succ}(n)$, then by (2.13.2) we get $\neg(0 = \text{succ}(n))$.

For the inductive step, let $x \equiv \text{succ}(n)$. If $y \equiv 0$, we use (2.13.2) again. Finally, if $y \equiv \text{succ}(m)$, the inductive hypothesis gives $(m = n) + \neg(m = n)$. In the first case, if $p : m = n$, then $\text{succ}(p) : \text{succ}(m) = \text{succ}(n)$. And in the second case, (2.13.3) yields $\neg(\text{succ}(m) = \text{succ}(n))$. \square

Although Hedberg's theorem appears rather special to sets (0-types), “Axiom K” generalizes naturally to n -types. Note that the ordinary Axiom K (as a property of a type X) states that for all $x : X$, the loop space $\Omega(X, x)$ (see Definition 2.1.8) is contractible. Since $\Omega(X, x)$ is always inhabited (by refl_x), this is equivalent to its being a mere proposition (a (-1) -type). Since $0 = (-1) + 1$, this suggests the following generalization.

Theorem 7.2.7. *For any $n \geq -1$, a type X is an $(n + 1)$ -type if and only if for all $x : X$, the type $\Omega(X, x)$ is an n -type.*

Before proving this, we prove an auxiliary lemma:

Lemma 7.2.8. *Given $n \geq -1$ and $X : \mathcal{U}$. If, given any inhabitant of X it follows that X is an n -type, then X is an n -type.*

Proof. Let $f : X \rightarrow \text{is-}n\text{-type}(X)$ be the given map. We need to show that for any $x, x' : X$, the type $x = x'$ is an $(n - 1)$ -type. But then $f(x)$ shows that X is an n -type, hence all its path spaces are $(n - 1)$ -types. \square

Proof of Theorem 7.2.7. The “only if” direction is obvious, since $\Omega(X, x) := (x =_X x)$. Conversely, in order to show that X is an $(n + 1)$ -type, we need to show that for any $x, x' : X$, the type $x = x'$ is an n -type. Following Lemma 7.2.8 it suffices to give a map

$$(x = x') \rightarrow \text{is-}n\text{-type}(x = x').$$

By path induction, it suffices to do this when $x \equiv x'$, in which case it follows from the assumption that $\Omega(X, x)$ is an n -type. \square

By induction and some slightly clever whiskering, we can obtain a generalization of the K property to $n > 0$.

Theorem 7.2.9. *For every $n \geq -1$, a type A is an n -type if and only if $\Omega^{n+1}(A, a)$ is contractible for all $a : A$.*

Proof. Recalling that $\Omega^0(A, a) = (A, a)$, the case $n = -1$ is Exercise 3.5. The case $n = 0$ is Theorem 7.2.1. Now we use induction; suppose the statement holds for $n : \mathbb{N}$. By Theorem 7.2.7, A is an $(n + 1)$ -type iff $\Omega(A, a)$ is an n -type for all $a : A$. By the inductive hypothesis, the latter is equivalent to saying that $\Omega^{n+1}(\Omega(A, a), p)$ is contractible for all $p : \Omega(A, a)$.

Since $\Omega^{n+2}(A, a) := \Omega^{n+1}(\Omega(A, a), \text{refl}_a)$, and $\Omega^{n+1} = \Omega^n \circ \Omega$, it will suffice to show that $\Omega(\Omega(A, a), p)$ is equal to $\Omega(\Omega(A, a), \text{refl}_a)$, in the type \mathcal{U}_\bullet of pointed types. For this, it suffices to give an equivalence

$$g : \Omega(\Omega(A, a), p) \simeq \Omega(\Omega(A, a), \text{refl}_a)$$

which carries the basepoint refl_p to the basepoint $\text{refl}_{\text{refl}_a}$. For $q : p = p$, define $g(q) : \text{refl}_a = \text{refl}_a$ to be the following composite:

$$\text{refl}_a = p \cdot p^{-1} \stackrel{q}{=} p \cdot p^{-1} = \text{refl}_a,$$

where the path labeled “ q ” is actually $\text{ap}_{\lambda r.r \cdot p^{-1}}(q)$. Then g is an equivalence because it is a composite of equivalences

$$(p = p) \xrightarrow{\text{ap}_{\lambda r.r \cdot p^{-1}}} (p \cdot p^{-1} = p \cdot p^{-1}) \xrightarrow{i \cdot - \cdot i^{-1}} (\text{refl}_a = \text{refl}_a).$$

using Example 2.4.8 and Theorem 2.11.1, where $i : \text{refl}_a = p \cdot p^{-1}$ is the canonical equality. And it is evident that $g(\text{refl}_p) = \text{refl}_{\text{refl}_a}$. \square

7.3 Truncations

In §3.7 we introduced the propositional truncation, which makes the “best approximation” of a type that is a mere proposition, i.e. a (-1) -type. In §6.9 we constructed this truncation as a higher inductive type, and gave one way to generalize it to a 0-truncation. We now explain a better generalization of this, which truncates any type into an n -type for any $n \geq -2$; in classical homotopy theory this would be called its n^{th} **Postnikov section**.

The idea is to make use of Theorem 7.2.9, which states that A is an n -type just when $\Omega^{n+1}(A, a)$ is contractible for all $a : A$, and Lemma 6.5.4, which implies that $\Omega^{n+1}(A, a) \simeq \text{Map}_*(\mathbb{S}^{n+1}, (A, a))$, where \mathbb{S}^{n+1} is equipped with some basepoint which we may as well call *base*. However, contractibility of $\text{Map}_*(\mathbb{S}^{n+1}, (A, a))$ is something that we can ensure directly by giving path constructors.

We will use the “hub and spoke” construction as in §6.7. Thus, for $n \geq -1$, we take $\|A\|_n$ to be the higher inductive type generated by:

- a function $|-|_n : A \rightarrow \|A\|_n$,
- for each $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$, a *hub* point $h(r) : \|A\|_n$, and
- for each $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ and each $x : \mathbb{S}^{n+1}$, a *spoke* path $s_r(x) : r(x) = h(r)$.

The existence of these constructors is now enough to show:

Lemma 7.3.1. $\|A\|_n$ is an n -type.

Proof. By Theorem 7.2.9, it suffices to show that $\Omega^{n+1}(\|A\|_n, b)$ is contractible for all $b : \|A\|_n$, which by Lemma 6.5.4 is equivalent to $\text{Map}_*(\mathbb{S}^{n+1}, (\|A\|_n, b))$. As center of contraction for the latter, we choose the function $c_b : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ which is constant at b , together with $\text{refl}_b : c_b(\text{base}) = b$.

Now, an arbitrary element of $\text{Map}_*(\mathbb{S}^{n+1}, (\|A\|_n, b))$ consists of a map $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ together with a path $p : r(\text{base}) = b$. By function extensionality, to show $r = c_b$ it suffices to give, for each $x : \mathbb{S}^{n+1}$, a path $r(x) = c_b(x) \equiv b$. We choose this to be the composite $s_r(x) \bullet s_r(\text{base})^{-1} \bullet p$, where $s_r(x)$ is the spoke at x .

Finally, we must show that when transported along this equality $r = c_b$, the path p becomes refl_b . By transport in path types, this means we need

$$(s_r(\text{base}) \bullet s_r(\text{base})^{-1} \bullet p)^{-1} \bullet p = \text{refl}_b.$$

But this is immediate from path operations. \square

(This construction fails for $n = -2$, but in that case we can simply define $\|A\|_{-2} := \mathbf{1}$ for all A . From now on we assume $n \geq -1$.)

To show the desired universal property of the n -truncation, we need the induction principle. We extract this from the constructors in the usual way; it says that given $P : \|A\|_n \rightarrow \mathcal{U}$ together with

- For each $a : A$, an element $g(a) : P(|a|_n)$,
- For each $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ and $r' : \prod_{(x:\mathbb{S}^{n+1})} P(r(x))$, an element $h'(r, r') : P(h(r))$,
- For each $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ and $r' : \prod_{(x:\mathbb{S}^{n+1})} P(r(x))$, and each $x : \mathbb{S}^{n+1}$, a dependent path $r'(x) =_{s_r(x)}^P h'(r, r')$,

there exists a section $f : \prod_{(x:\|A\|_n)} P(x)$ with $f(|a|_n) \equiv g(a)$ for all $a : A$. To make this more useful, we reformulate it as follows.

Theorem 7.3.2. *For any type family $P : \|A\|_n \rightarrow \mathcal{U}$ such that each $P(x)$ is an n -type, and any function $g : \prod_{(a:A)} P(|a|_n)$, there exists a section $f : \prod_{(x:\|A\|_n)} P(x)$ such that $f(|a|_n) \equiv g(a)$ for all $a : A$.*

Proof. It will suffice to construct the second and third data listed above, since g has exactly the type of the first datum. Given $r : \mathbb{S}^{n+1} \rightarrow \|A\|_n$ and $r' : \prod_{(x:\mathbb{S}^{n+1})} P(r(x))$, we have $h(r) : \|A\|_n$ and $s_r : \prod_{(x:\mathbb{S}^{n+1})} (r(x) = h(r))$. Define $t : \mathbb{S}^{n+1} \rightarrow P(h(r))$ by $t(x) := s_r(x)_*(r'(x))$. Then since $P(h(r))$ is n -truncated, there exists a point $u : P(h(r))$ and a contraction $v : \prod_{(x:\mathbb{S}^{n+1})} (t(x) = u)$. Define $h'(r, r') := u$, giving the second datum. Then (recalling the definition of dependent paths), v has exactly the type required of the third datum. \square

In particular, if E is some n -type, we can consider the constant family of types equal to E for every point of A . Thus, every map $f : A \rightarrow E$ can be extended to a map $\text{ext}(f) : \|A\|_n \rightarrow E$ defined by $\text{ext}(f)(|a|_n) := f(a)$; this is the *recursion principle* for $\|A\|_n$.

The induction principle also implies a uniqueness principle for functions of this form. Namely, if E is an n -type and $g, g' : \|A\|_n \rightarrow E$ are such that $g(|a|_n) = g'(|a|_n)$ for every $a : A$, then $g(x) = g'(x)$ for all $x : \|A\|_n$, since the type $g(x) = g'(x)$ is an n -type. Thus, $g = g'$. (In fact, this uniqueness principle holds more generally when E is an $(n+1)$ -type.) This yields the following universal property.

Lemma 7.3.3 (Universal property of truncations). *Let $n \geq -2$, $A : \mathcal{U}$ and $B : n\text{-Type}$. The following map is an equivalence:*

$$\begin{cases} (\|A\|_n \rightarrow B) & \longrightarrow (A \rightarrow B) \\ g & \longmapsto g \circ |-|_n \end{cases}$$

Proof. Given that B is n -truncated, any $f : A \rightarrow B$ can be extended to a map $\text{ext}(f) : \|A\|_n \rightarrow B$. The map $\text{ext}(f) \circ |-|_n$ is equal to f , because for every $a : A$ we have $\text{ext}(f)(|a|_n) = f(a)$ by definition. And the map $\text{ext}(g \circ |-|_n)$ is equal to g , because they both send $|a|_n$ to $g(|a|_n)$. \square

In categorical language, this says that the n -types form a *reflective subcategory* of the category of types. (To state this fully precisely, one ought to use the language of $(\infty, 1)$ -categories.) In particular, this implies that the n -truncation is functorial: given $f : A \rightarrow B$, applying the recursion principle to the composite $A \xrightarrow{f} B \rightarrow \|B\|_n$ yields a map $\|f\|_n : \|A\|_n \rightarrow \|B\|_n$. By definition, we have a homotopy

$$\text{nat}_n^f : \prod_{a:A} \|f\|_n(|a|_n) = |f(a)|_n, \quad (7.3.4)$$

expressing *naturality* of the maps $|-|_n$.

Uniqueness implies functoriality laws such as $\|g \circ f\|_n = \|g\|_n \circ \|f\|_n$ and $\|\text{id}_A\|_n = \text{id}_{\|A\|_n}$, with attendant coherence laws. We also have higher functoriality, for instance:

Lemma 7.3.5. *Given $f, g : A \rightarrow B$ and a homotopy $h : f \sim g$, there is an induced homotopy $\|h\|_n : \|f\|_n \sim \|g\|_n$ such that the composite*

$$|f(a)|_n \xrightarrow{\text{nat}_n^f(a)^{-1}} \|f\|_n(|a|_n) \xrightarrow{\|h\|_n(|a|_n)} \|g\|_n(|a|_n) \xrightarrow{\text{nat}_n^g(a)} |g(a)|_n \quad (7.3.6)$$

is equal to $\text{ap}_{|-|_n}(h(a))$.

Proof. First, we indeed have a homotopy with components $\text{ap}_{|-|_n}(h(a)) : |f(a)|_n = |g(a)|_n$. Composing on either sides with the paths $|f(a)|_n = \|f\|_n(|a|_n)$ and $|g(a)|_n = \|g\|_n(|a|_n)$, which arise from the definitions of $\|f\|_n$ and $\|g\|_n$, we obtain a homotopy $(\|f\|_n \circ |-|_n) \sim (\|g\|_n \circ |-|_n)$, and hence an equality by function extensionality. But since $(- \circ |-|_n)$ is an equivalence, there must be a path $\|f\|_n = \|g\|_n$ inducing it, and the coherence laws for function extensionality imply (7.3.6). \square

The following observation about reflective subcategories is also standard.

Corollary 7.3.7. *A type A is an n -type if and only if $|-|_n : A \rightarrow \|A\|_n$ is an equivalence.*

Proof. “If” follows from closure of n -types under equivalence. On the other hand, if A is an n -type, we can define $\text{ext}(\text{id}_A) : \|A\|_n \rightarrow A$. Then we have $\text{ext}(\text{id}_A) \circ |-|_n = \text{id}_A : A \rightarrow A$ by definition. In order to prove that $|-|_n \circ \text{ext}(\text{id}_A) = \text{id}_{\|A\|_n}$, we only need to prove that $|-|_n \circ \text{ext}(\text{id}_A) \circ |-|_n = \text{id}_{\|A\|_n} \circ |-|_n$. This is again true:

$$\begin{array}{ccc} A & \xrightarrow{|-|_n} & \|A\|_n \\ & \searrow \text{id}_A & \downarrow \text{ext}(\text{id}_A) \\ & A & \\ & \downarrow |-|_n & \swarrow \text{id}_{\|A\|_n} \\ & \|A\|_n & \end{array}$$

\square

The category of n -types also has some special properties not possessed by all reflective subcategories. For instance, the reflector $\|-|_n$ preserves finite products.

Theorem 7.3.8. *For any types A and B , the induced map $\|A \times B\|_n \rightarrow \|A\|_n \times \|B\|_n$ is an equivalence.*

Proof. It suffices to show that $\|A\|_n \times \|B\|_n$ has the same universal property as $\|A \times B\|_n$. Thus, let C be an n -type; we have

$$\begin{aligned} (\|A\|_n \times \|B\|_n \rightarrow C) &= (\|A\|_n \rightarrow (\|B\|_n \rightarrow C)) \\ &= (\|A\|_n \rightarrow (B \rightarrow C)) \\ &= (A \rightarrow (B \rightarrow C)) \\ &= (A \times B \rightarrow C) \end{aligned}$$

using the universal properties of $\|B\|_n$ and $\|A\|_n$, along with the fact that $B \rightarrow C$ is an n -type since C is. It is straightforward to verify that this equivalence is given by composing with $|-|_n \times |-|_n$, as needed. \square

The following related fact about dependent sums is often useful.

Theorem 7.3.9. *Let $P : A \rightarrow \mathcal{U}$ be a family of types. Then there is an equivalence*

$$\left\| \sum_{x:A} \|P(x)\|_n \right\|_n \simeq \left\| \sum_{x:A} P(x) \right\|_n.$$

Proof. We use the induction principle of n -truncation several times to construct functions

$$\begin{aligned}\varphi : \left\| \sum_{x:A} \|P(x)\|_n \right\|_n &\rightarrow \left\| \sum_{x:A} P(x) \right\|_n \\ \psi : \left\| \sum_{x:A} P(x) \right\|_n &\rightarrow \left\| \sum_{x:A} \|P(x)\|_n \right\|_n\end{aligned}$$

and homotopies $H : \varphi \circ \psi \sim \text{id}$ and $K : \psi \circ \varphi \sim \text{id}$ exhibiting them as quasi-inverses. We define φ by setting $\varphi(|(x, |u|_n)|_n) := |(x, u)|_n$. We define ψ by setting $\psi(|(x, u)|_n) := |(x, |u|_n)|_n$. Then we define $H(|(x, u)|_n) := \text{refl}_{|(x, u)|_n}$ and $K(|(x, |u|_n)|_n) := \text{refl}_{|(x, |u|_n)|_n}$. \square

Corollary 7.3.10. *If A is an n -type and $P : A \rightarrow \mathcal{U}$ is any type family, then*

$$\sum_{a:A} \|P(a)\|_n \simeq \left\| \sum_{a:A} P(a) \right\|_n$$

Proof. If A is an n -type, then the left-hand type above is already an n -type, hence equivalent to its n -truncation; thus this follows from Theorem 7.3.9. \square

We can characterize the path spaces of a truncation using the same method that we used in §§2.12 and 2.13 for coproducts and natural numbers (and which we will use in ?? to calculate homotopy groups). Unsurprisingly, the path spaces in the $(n+1)$ -truncation of A are the n -truncations of the path spaces of A . Indeed, for any $x, y : A$ there is a canonical map

$$f : \|x =_A y\|_n \rightarrow \left(|x|_{n+1} =_{\|A\|_{n+1}} |y|_{n+1} \right) \quad (7.3.11)$$

defined by

$$f(|p|_n) := \text{ap}_{|-|_{n+1}}(p).$$

This definition uses the recursion principle for $\|-|_n$, which is correct because $\|A\|_{n+1}$ is $(n+1)$ -truncated, so that the codomain of f is n -truncated.

Theorem 7.3.12. *For any A and $x, y : A$ and $n \geq -2$, the map (7.3.11) is an equivalence; thus we have*

$$\|x =_A y\|_n \simeq \left(|x|_{n+1} =_{\|A\|_{n+1}} |y|_{n+1} \right).$$

Proof. The proof is a simple application of the encode-decode method: As in previous situations, we cannot directly define a quasi-inverse to the map (7.3.11) because there is no way to induct on an equality between $|x|_{n+1}$ and $|y|_{n+1}$. Thus, instead we generalize its type, in order to have general elements of the type $\|A\|_{n+1}$ instead of $|x|_{n+1}$ and $|y|_{n+1}$. Define $P : \|A\|_{n+1} \rightarrow \|A\|_{n+1} \rightarrow n\text{-Type}$ by

$$P(|x|_{n+1}, |y|_{n+1}) := \|x =_A y\|_n$$

This definition is correct because $\|x =_A y\|_n$ is n -truncated, and $n\text{-Type}$ is $(n+1)$ -truncated by Theorem 7.1.11. Now for every $u, v : \|A\|_{n+1}$, there is a map

$$\text{decode} : P(u, v) \rightarrow (u =_{\|A\|_{n+1}} v)$$

defined for $u = |x|_{n+1}$ and $v = |y|_{n+1}$ and $p : x = y$ by

$$\text{decode}(|p|_n) := \text{ap}_{|-|_{n+1}}(p).$$

Since the codomain of decode is n -truncated, it suffices to define it only for u and v of this form, and then it's just the same definition as before. We also define a function

$$r : \prod_{u:\|A\|_{n+1}} P(u, u)$$

by induction on u , where $r(|x|_{n+1}) := |\text{refl}_x|_n$.

Now we can define an inverse map

$$\text{encode} : (u =_{\|A\|_{n+1}} v) \rightarrow P(u, v)$$

by

$$\text{encode}(p) := \text{transport}^{v \mapsto P(u, v)}(p, r(u)).$$

To show that the composite

$$(u =_{\|A\|_{n+1}} v) \xrightarrow{\text{encode}} P(u, v) \xrightarrow{\text{decode}} (u =_{\|A\|_{n+1}} v)$$

is the identity function, by path induction it suffices to check it for $\text{refl}_u : u = u$, in which case what we need to know is that $\text{decode}(r(u)) = \text{refl}_u$. But since this is an $(n-1)$ -type, hence also an $(n+1)$ -type, we may assume $u \equiv |x|_{n+1}$, in which case it follows by definition of r and decode . Finally, to show that

$$P(u, v) \xrightarrow{\text{decode}} (u =_{\|A\|_{n+1}} v) \xrightarrow{\text{encode}} P(u, v)$$

is the identity function, since this goal is again an $(n-1)$ -type, we may assume that $u = |x|_{n+1}$ and $v = |y|_{n+1}$ and that we are considering $|p|_n : P(|x|_{n+1}, |y|_{n+1})$ for some $p : x = y$. Then we have

$$\begin{aligned} \text{encode}(\text{decode}(|p|_n)) &= \text{encode}(\text{ap}_{|-|_{n+1}}(p)) \\ &= \text{transport}^{v \mapsto P(|x|_{n+1}, v)}(\text{ap}_{|-|_{n+1}}(p), |\text{refl}_x|_n) \\ &= \text{transport}^{y \mapsto \|x=y\|_n}(p, |\text{refl}_x|_n) \\ &= \left| \text{transport}^{y \mapsto (x=y)}(p, \text{refl}_x) \right|_n \\ &= |p|_n, \end{aligned}$$

using Lemmas 2.3.10 and 2.3.11. (Alternatively, we could do path induction on p ; the desired equality would then hold judgmentally.) This completes the proof that decode and encode are quasi-inverses. The stated result is then the special case where $u = |x|_{n+1}$ and $v = |y|_{n+1}$. \square

Corollary 7.3.13. *Let $n \geq -2$ and (A, a) be a pointed type. Then*

$$\|\Omega(A, a)\|_n = \Omega(\|(A, a)\|_{n+1})$$

Proof. This is a special case of the previous lemma where $x = y = a$. \square

Corollary 7.3.14. *Let $n \geq -2$ and $k \geq 0$ and (A, a) a pointed type. Then*

$$\|\Omega^k(A, a)\|_n = \Omega^k(\|(A, a)\|_{n+k}).$$

Proof. By induction on k , using the recursive definition of Ω^k . \square

We also observe that “truncations are cumulative”: if we truncate to an n -type and then to a k -type with $k \leq n$, then we might as well have truncated directly to a k -type.

Lemma 7.3.15. *Let $k, n \geq -2$ with $k \leq n$ and $A : \mathcal{U}$. Then $\|\|A\|_n\|_k = \|A\|_k$.*

Proof. We define two maps $f : \|\|A\|_n\|_k \rightarrow \|A\|_k$ and $g : \|A\|_k \rightarrow \|\|A\|_n\|_k$ by

$$f(\|a\|_n|_k) := |a|_k \quad \text{and} \quad g(|a|_k) := \|a\|_n|_k.$$

The map f is well-defined because $\|A\|_k$ is k -truncated and also n -truncated (because $k \leq n$), and the map g is well-defined because $\|\|A\|_n\|_k$ is k -truncated.

The composition $f \circ g : \|A\|_k \rightarrow \|A\|_k$ satisfies $(f \circ g)(|a|_k) = |a|_k$, hence $f \circ g = \text{id}_{\|A\|_k}$. Similarly, we have $(g \circ f)(\|a\|_n|_k) = \|a\|_n|_k$ and hence $g \circ f = \text{id}_{\|\|A\|_n\|_k}$. \square

7.4 Colimits of n -types

Recall that in §6.8, we used higher inductive types to define pushouts of types, and proved their universal property. In general, a (homotopy) colimit of n -types may no longer be an n -type (for an extreme counterexample, see Exercise 7.2). However, if we n -truncate it, we obtain an n -type which satisfies the correct universal property with respect to other n -types.

In this section we prove this for pushouts, which are the most important and nontrivial case of colimits. Recall the following definitions from §6.8.

Definition 7.4.1. A **span** is a 5-tuple $\mathcal{D} = (A, B, C, f, g)$ with $f : C \rightarrow A$ and $g : C \rightarrow B$.

$$\mathcal{D} = \begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & & \\ A & & \end{array}$$

Definition 7.4.2. Given a span $\mathcal{D} = (A, B, C, f, g)$ and a type D , a **cocone under \mathcal{D} with base D** is a triple (i, j, h) with $i : A \rightarrow D$, $j : B \rightarrow D$ and $h : \prod_{(c:C)} i(f(c)) = j(g(c))$:

$$\begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & \nearrow h & \downarrow j \\ A & \xrightarrow{i} & D \end{array}$$

We denote by $\text{cocone}_{\mathcal{D}}(D)$ the type of all such cocones.

The type of cocones is (covariantly) functorial. For instance, given D, E and a map $t : D \rightarrow E$, there is a map

$$\left\{ \begin{array}{rcl} \text{cocone}_{\mathcal{D}}(D) & \longrightarrow & \text{cocone}_{\mathcal{D}}(E) \\ c & \longmapsto & t \circ c \end{array} \right.$$

defined by:

$$t \circ (i, j, h) = (t \circ i, t \circ j, \text{ap}_t \circ h).$$

And given D, E, F , functions $t : D \rightarrow E$, $u : E \rightarrow F$ and $c : \text{cocone}_{\mathcal{D}}(D)$, we have

$$\text{id}_D \circ c = c \tag{7.4.3}$$

$$(u \circ t) \circ c = u \circ (t \circ c). \tag{7.4.4}$$

Definition 7.4.5. Given a span \mathcal{D} of n -types, an n -type D , and a cocone $c : \text{cocone}_{\mathcal{D}}(D)$, the pair (D, c) is said to be a **pushout of \mathcal{D} in n -types** if for every n -type E , the map

$$\begin{cases} (D \rightarrow E) & \longrightarrow \text{cocone}_{\mathcal{D}}(E) \\ t & \longmapsto t \circ c \end{cases}$$

is an equivalence.

In order to construct pushouts of n -types, we need to explain how to reflect spans and cocones.

Definition 7.4.6. Let

$$\mathcal{D} = \begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & & \\ A & & \end{array}$$

be a span. We denote by $\|\mathcal{D}\|_n$ the following span of n -types:

$$\|\mathcal{D}\|_n := \begin{array}{ccc} \|C\|_n & \xrightarrow{\|g\|_n} & \|B\|_n \\ \|f\|_n \downarrow & & \\ \|A\|_n & & \end{array}$$

Definition 7.4.7. Let $D : \mathcal{U}$ and $c = (i, j, h) : \text{cocone}_{\mathcal{D}}(D)$. We define

$$\|c\|_n = (\|i\|_n, \|j\|_n, k) : \text{cocone}_{\|\mathcal{D}\|_n}(\|D\|_n)$$

where k is the composite homotopy

$$\|i\|_n \circ \|f\|_n \sim \|i \circ f\|_n \sim \|j \circ g\|_n \sim \|j\|_n \circ \|g\|_n$$

using Lemma 7.3.5 and the functoriality of $\|- \|_n$.

We now observe that the maps from each type to its n -truncation assemble into a map of spans, in the following sense.

Definition 7.4.8. Let

$$\mathcal{D} = \begin{array}{ccc} C & \xrightarrow{g} & B \\ f \downarrow & & \\ A & & \end{array} \quad \text{and} \quad \mathcal{D}' = \begin{array}{ccc} C' & \xrightarrow{g'} & B' \\ f' \downarrow & & \\ A' & & \end{array}$$

be spans. A **map of spans** $\mathcal{D} \rightarrow \mathcal{D}'$ consists of functions $\alpha : A \rightarrow A'$, $\beta : B \rightarrow B'$, and $\gamma : C \rightarrow C'$ and homotopies $\phi : \alpha \circ f \sim f' \circ \gamma$ and $\psi : \beta \circ g \sim g' \circ \gamma$.

Thus, for any span \mathcal{D} , we have a map of spans $|-|_{\mathcal{D}}^n : \mathcal{D} \rightarrow \|\mathcal{D}\|_n$ consisting of $|-|_n^A$, $|-|_n^B$, $|-|_n^C$, and the naturality homotopies nat_n^f and nat_n^g from (7.3.4).

We also need to know that maps of spans behave functorially. Namely, if $(\alpha, \beta, \gamma, \phi, \psi) : \mathcal{D} \rightarrow \mathcal{D}'$ is a map of spans and D any type, then we have

$$\begin{cases} \text{cocone}_{\mathcal{D}'}(D) & \longrightarrow \text{cocone}_{\mathcal{D}}(D) \\ (i, j, h) & \longmapsto (i \circ \alpha, j \circ \beta, k) \end{cases}$$

where $k : \prod_{(z:C)} i(\alpha(f(z))) = j(\beta(g(z)))$ is the composite

$$i(\alpha(f(z))) \xrightarrow{\text{ap}_i(\phi)} i(f'(\gamma(z))) \xrightarrow{h(\gamma(z))} j(g'(\gamma(z))) \xrightarrow{\text{ap}_j(\psi)} j(\beta(g(z))). \quad (7.4.9)$$

We denote this cocone by $(i, j, h) \circ (\alpha, \beta, \gamma, \phi, \psi)$. Moreover, this functorial action commutes with the other functoriality of cocones:

Lemma 7.4.10. *Given $(\alpha, \beta, \gamma, \phi, \psi) : \mathcal{D} \rightarrow \mathcal{D}'$ and $t : D \rightarrow E$, the following diagram commutes:*

$$\begin{array}{ccc} \text{cocone}_{\mathcal{D}'}(D) & \xrightarrow{t \circ -} & \text{cocone}_{\mathcal{D}'}(E) \\ \downarrow & & \downarrow \\ \text{cocone}_{\mathcal{D}}(D) & \xrightarrow{t \circ -} & \text{cocone}_{\mathcal{D}}(E) \end{array}$$

Proof. Given $(i, j, h) : \text{cocone}_{\mathcal{D}'}(D)$, note that both composites yield a cocone whose first two components are $t \circ i \circ \alpha$ and $t \circ j \circ \beta$. Thus, it remains to verify that the homotopies agree. For the top-right composite, the homotopy is (7.4.9) with (i, j, h) replaced by $(t \circ i, t \circ j, \text{ap}_t \circ h)$:

$$t i \alpha f z \xrightarrow{\text{ap}_{t \circ i}(\phi)} t i f' \gamma z \xrightarrow{\text{ap}_t(h(\gamma(z)))} t j g' \gamma z \xrightarrow{\text{ap}_{t \circ j}(\psi)} t j \beta g z$$

(For brevity, we are omitting the parentheses around the arguments of functions.) On the other hand, for the left-bottom composite, the homotopy is ap_t applied to (7.4.9). Since ap respects path-concatenation, this is equal to

$$t i \alpha f z \xrightarrow{\text{ap}_t(\text{ap}_i(\phi))} t i f' \gamma z \xrightarrow{\text{ap}_t(h(\gamma(z)))} t j g' \gamma z \xrightarrow{\text{ap}_t(\text{ap}_j(\psi))} t j \beta g z.$$

But $\text{ap}_t \circ \text{ap}_i = \text{ap}_{t \circ i}$ and similarly for j , so these two homotopies are equal. \square

Finally, note that since we defined $\|c\|_n : \text{cocone}_{\|\mathcal{D}\|_n}(\|D\|_n)$ using Lemma 7.3.5, the additional condition (7.3.6) implies

$$|-|_n^D \circ c = \|c\|_n \circ |-|_n^{\mathcal{D}}. \quad (7.4.11)$$

for any $c : \text{cocone}_{\mathcal{D}}(D)$. Now we can prove our desired theorem.

Theorem 7.4.12. *Let \mathcal{D} be a span and (D, c) its pushout. Then $(\|D\|_n, \|c\|_n)$ is a pushout of $\|\mathcal{D}\|_n$ in n -types.*

Proof. Let E be an n -type, and consider the following diagram:

$$\begin{array}{ccc} (\|D\|_n \rightarrow E) & \xrightarrow{-\circ|-\|_n^D} & (D \rightarrow E) \\ \downarrow -\circ\|c\|_n & & \downarrow -\circ c \\ \text{cocone}_{\|\mathcal{D}\|_n}(E) & \xrightarrow{-\circ|-\|_n^{\mathcal{D}}} & \text{cocone}_{\mathcal{D}}(E) \\ \uparrow \ell_1 & & \uparrow \ell_2 \\ (\|A\|_n \rightarrow E) \times_{(\|C\|_n \rightarrow E)} (\|B\|_n \rightarrow E) & \longrightarrow & (A \rightarrow E) \times_{(C \rightarrow E)} (B \rightarrow E) \end{array}$$

The upper horizontal arrow is an equivalence since E is an n -type, while $- \circ c$ is an equivalence since c is a pushout cocone. Thus, by the 2-out-of-3 property, to show that $- \circ \|c\|_n$ is an equivalence, it will suffice to show that the upper square commutes and that the middle horizontal arrow is an equivalence. To see that the upper square commutes, let $t : \|D\|_n \rightarrow E$; then

$$\begin{aligned} (t \circ \|c\|_n) \circ |-|_n^{\mathcal{D}} &= t \circ (\|c\|_n \circ |-|_n^{\mathcal{D}}) && \text{(by Lemma 7.4.10)} \\ &= t \circ (|-|_n^D \circ c) && \text{(by (7.4.11))} \\ &= (t \circ |-|_n^D) \circ c. && \text{(by (7.4.4))} \end{aligned}$$

To show that the middle horizontal arrow is an equivalence, consider the lower square. The two lower vertical arrows are simply applications of `happly`:

$$\begin{aligned} \ell_1(i, j, p) &:= (i, j, \text{happly}(p)) \\ \ell_2(i, j, p) &:= (i, j, \text{happly}(p)) \end{aligned}$$

and hence are equivalences by function extensionality. The lowest horizontal arrow is defined by

$$(i, j, p) \mapsto (i \circ |-|_n^A, j \circ |-|_n^B, q)$$

where q is the composite

$$\begin{aligned} i \circ |-|_n^A \circ f &= i \circ \|f\|_n \circ |-|_n^C && \text{(by } \text{funext}(\lambda z. \text{ap}_i(\text{nat}_n^f(z)))\text{)} \\ &= j \circ \|g\|_n \circ |-|_n^C && \text{(by } \text{ap}_{- \circ |-|_n^C}(p)\text{)} \\ &= j \circ |-|_n^B \circ g. && \text{(by } \text{funext}(\lambda z. \text{ap}_j(\text{nat}_n^g(z)))\text{)} \end{aligned}$$

This is an equivalence, because it is induced by an equivalence of cospans. Thus, by 2-out-of-3, it will suffice to show that the lower square commutes. But the two composites around the lower square agree definitionally on the first two components, so it suffices to show that for (i, j, p) in the lower left corner and $z : C$, the path

$$\text{happly}(q, z) : i(|f(z)|_n) = j(|g(z)|_n)$$

(with q as above) is equal to the composite

$$\begin{aligned} i(|f(z)|_n) &= i(\|f\|_n(|z|_n)) && \text{(by } \text{ap}_i(\text{nat}_n^f(z))\text{)} \\ &= j(\|g\|_n(|z|_n)) && \text{(by } \text{happly}(p, |z|_n)\text{)} \\ &= j(|g(z)|_n). && \text{(by } \text{ap}_j(\text{nat}_n^g(z))\text{)} \end{aligned}$$

However, since `happly` is functorial, it suffices to check equality for the three component paths:

$$\begin{aligned} \text{happly}(\text{funext}(\lambda z. \text{ap}_i(\text{nat}_n^f(z))), z) &= \text{ap}_i(\text{nat}_n^f(z)) \\ \text{happly}(\text{ap}_{- \circ |-|_n^C}(p), z) &= \text{happly}(p, |z|_n) \\ \text{happly}(\text{funext}(\lambda z. \text{ap}_j(\text{nat}_n^g(z))), z) &= \text{ap}_j(\text{nat}_n^g(z)). \end{aligned}$$

The first and third of these are just the fact that `happly` is quasi-inverse to `funext`, while the second is an easy general lemma about `happly` and precomposition. \square

7.5 Connectedness

An n -type is one that has no interesting information above dimension n . By contrast, an n -connected type is one that has no interesting information *below* dimension n . It turns out to be natural to study a more general notion for functions as well.

Definition 7.5.1. A function $f : A \rightarrow B$ is said to be **n -connected** if for all $b : B$, the type $\|\text{fib}_f(b)\|_n$ is contractible:

$$\text{conn}_n(f) : \equiv \prod_{b:B} \text{isContr}(\|\text{fib}_f(b)\|_n).$$

A type A is said to be **n -connected** if the unique function $A \rightarrow \mathbf{1}$ is n -connected, i.e. if $\|A\|_n$ is contractible.

Thus, a function $f : A \rightarrow B$ is n -connected if and only if $\text{fib}_f(b)$ is n -connected for every $b : B$. Of course, every function is (-2) -connected. At the next level, we have:

Lemma 7.5.2. *A function f is (-1) -connected if and only if it is surjective in the sense of §4.6.*

Proof. We defined f to be surjective if $\|\text{fib}_f(b)\|_{-1}$ is inhabited for all b . But since it is a mere proposition, inhabitation is equivalent to contractibility. \square

Thus, n -connectedness of a function for $n \geq 0$ can be thought of as a strong form of surjectivity. Category-theoretically, (-1) -connectedness corresponds to essential surjectivity on objects, while n -connectedness corresponds to essential surjectivity on k -morphisms for $k \leq n + 1$.

Lemma 7.5.2 also implies that a type A is (-1) -connected if and only if it is merely inhabited. When a type is 0-connected we may simply say that it is **connected**, and when it is 1-connected we say it is **simply connected**.

Remark 7.5.3. While our notion of n -connectedness for types agrees with the standard notion in homotopy theory, our notion of n -connectedness for *functions* is off by one from a common indexing in classical homotopy theory. Whereas we say a function f is n -connected if all its fibers are n -connected, some classical homotopy theorists would call such a function $(n + 1)$ -connected. (This is due to a historical focus on *cofibers* rather than fibers.)

We now observe a few closure properties of connected maps.

Lemma 7.5.4. *Suppose that g is a retract of a n -connected function f . Then g is n -connected.*

Proof. This is a direct consequence of Lemma 4.7.3. \square

Corollary 7.5.5. *If g is homotopic to a n -connected function f , then g is n -connected.*

Lemma 7.5.6. *Suppose that $f : A \rightarrow B$ is n -connected. Then $g : B \rightarrow C$ is n -connected if and only if $g \circ f$ is n -connected.*

Proof. For any $c : C$, we have

$$\begin{aligned} \|\text{fib}_{g \circ f}(c)\|_n &\simeq \left\| \sum_{w:\text{fib}_g(c)} \text{fib}_f(\text{pr}_1 w) \right\|_n && \text{(by Exercise 4.4)} \\ &\simeq \left\| \sum_{w:\text{fib}_g(c)} \|\text{fib}_f(\text{pr}_1 w)\|_n \right\|_n && \text{(by Theorem 7.3.9)} \\ &\simeq \|\text{fib}_g(c)\|_n. && \text{(since } \|\text{fib}_f(\text{pr}_1 w)\|_n \text{ is contractible)} \end{aligned}$$

It follows that $\|\text{fib}_g(c)\|_n$ is contractible if and only if $\|\text{fib}_{g \circ f}(c)\|_n$ is contractible. \square

Importantly, n -connected functions can be equivalently characterized as those which satisfy an “induction principle” with respect to n -types. This idea will lead directly into our proof of the Freudenthal suspension theorem in ??.

Lemma 7.5.7. *For $f : A \rightarrow B$ and $P : B \rightarrow \mathcal{U}$, consider the following function:*

$$\lambda s. s \circ f : \left(\prod_{b:B} P(b) \right) \rightarrow \left(\prod_{a:A} P(f(a)) \right).$$

For a fixed f and $n \geq -2$, the following are equivalent.

- (i) f is n -connected.
- (ii) For every $P : B \rightarrow n\text{-Type}$, the map $\lambda s. s \circ f$ is an equivalence.
- (iii) For every $P : B \rightarrow n\text{-Type}$, the map $\lambda s. s \circ f$ has a section.

Proof. Suppose that f is n -connected and let $P : B \rightarrow n\text{-Type}$. Then we have the equivalences

$$\begin{aligned} \prod_{b:B} P(b) &\simeq \prod_{b:B} \left(\|\text{fib}_f(b)\|_n \rightarrow P(b) \right) && \text{(since } \|\text{fib}_f(b)\|_n \text{ is contractible)} \\ &\simeq \prod_{b:B} \left(\text{fib}_f(b) \rightarrow P(b) \right) && \text{(since } P(b) \text{ is an } n\text{-type)} \\ &\simeq \prod_{(b:B)} \prod_{(a:A)} \prod_{(p:f(a)=b)} P(b) && \text{(by the left universal property of } \Sigma\text{-types)} \\ &\simeq \prod_{a:A} P(f(a)). && \text{(by the left universal property of path types)} \end{aligned}$$

We omit the proof that this equivalence is indeed given by $\lambda s. s \circ f$. Thus, (i) \Rightarrow (ii), and clearly (ii) \Rightarrow (iii). To show (iii) \Rightarrow (i), consider the type family

$$P(b) := \|\text{fib}_f(b)\|_n.$$

Then (iii) yields a map $c : \prod_{(b:B)} \|\text{fib}_f(b)\|_n$ with $c(f(a)) = |(a, \text{refl}_{f(a)})|_n$. To show that each $\|\text{fib}_f(b)\|_n$ is contractible, we will find a function of type

$$\prod_{(b:B)} \prod_{(w:\|\text{fib}_f(b)\|_n)} w = c(b).$$

By Theorem 7.3.2, for this it suffices to find a function of type

$$\prod_{(b:B)} \prod_{(a:A)} \prod_{(p:f(a)=b)} |(a, p)|_n = c(b).$$

But by rearranging variables and path induction, this is equivalent to the type

$$\prod_{a:A} |(a, \text{refl}_{f(a)})|_n = c(f(a)).$$

This property holds by our choice of $c(f(a))$. □

Corollary 7.5.8. *For any A , the canonical function $|-\|_n : A \rightarrow \|A\|_n$ is n -connected.*

Proof. By Theorem 7.3.2 and the associated uniqueness principle, the condition of Lemma 7.5.7 holds. □

For instance, when $n = -1$, Corollary 7.5.8 says that the map $A \rightarrow \|A\|$ from a type to its propositional truncation is surjective.

Corollary 7.5.9. *A type A is n -connected if and only if the map*

$$\lambda b. \lambda a. b : B \rightarrow (A \rightarrow B)$$

is an equivalence for every n -type B . In other words, “every map from A to an n -type is constant”.

Proof. By Lemma 7.5.7 applied to a function with codomain $\mathbf{1}$. \square

Lemma 7.5.10. *Let B be an n -type and let $f : A \rightarrow B$ be a function. Then the induced function $g : \|A\|_n \rightarrow B$ is an equivalence if and only if f is n -connected.*

Proof. By Corollary 7.5.8, $|-\|_n$ is n -connected. Thus, since $f = g \circ |-\|_n$, by Lemma 7.5.6 f is n -connected if and only if g is n -connected. But since g is a function between n -types, its fibers are also n -types. Thus, g is n -connected if and only if it is an equivalence. \square

We can also characterize connected pointed types in terms of connectivity of the inclusion of their basepoint.

Lemma 7.5.11. *Let A be a type and $a_0 : \mathbf{1} \rightarrow A$ a basepoint, with $n \geq -1$. Then A is n -connected if and only if the map a_0 is $(n-1)$ -connected.*

Proof. First suppose $a_0 : \mathbf{1} \rightarrow A$ is $(n-1)$ -connected and let B be an n -type; we will use Corollary 7.5.9. The map $\lambda b. \lambda a. b : B \rightarrow (A \rightarrow B)$ has a retraction given by $f \mapsto f(a_0)$, so it suffices to show it also has a section, i.e. that for any $f : A \rightarrow B$ there is $b : B$ such that $f = \lambda a. b$. We choose $b := f(a_0)$. Define $P : A \rightarrow \mathcal{U}$ by $P(a) := (f(a) = f(a_0))$. Then P is a family of $(n-1)$ -types and we have $P(a_0)$; hence we have $\prod_{(a:A)} P(a)$ since $a_0 : \mathbf{1} \rightarrow A$ is $(n-1)$ -connected. Thus, $f = \lambda a. f(a_0)$ as desired.

Now suppose A is n -connected, and let $P : A \rightarrow (n-1)\text{-Type}$ and $u : P(a_0)$ be given. By Lemma 7.5.7, it will suffice to construct $f : \prod_{(a:A)} P(a)$ such that $f(a_0) = u$. Now $(n-1)\text{-Type}$ is an n -type and A is n -connected, so by Corollary 7.5.9, there is an n -type B such that $P = \lambda a. B$. Hence, we have a family of equivalences $g : \prod_{(a:A)} (P(a) \simeq B)$. Define $f(a) := g_a^{-1}(g_{a_0}(u))$; then $f : \prod_{(a:A)} P(a)$ and $f(a_0) = u$ as desired. \square

In particular, a pointed type (A, a_0) is 0-connected if and only if $a_0 : \mathbf{1} \rightarrow A$ is surjective, which is to say $\prod_{(x:A)} \|x = a_0\|$. For a similar result in the not-necessarily-pointed case, see Exercise 7.6.

A useful variation on Lemma 7.5.6 is:

Lemma 7.5.12. *Let $f : A \rightarrow B$ be a function and $P : A \rightarrow \mathcal{U}$ and $Q : B \rightarrow \mathcal{U}$ be type families. Suppose that $g : \prod_{(a:A)} P(a) \rightarrow Q(f(a))$ is a fiberwise n -connected family of functions, i.e. each function $g_a : P(a) \rightarrow Q(f(a))$ is n -connected. If f is also n -connected, then so is the function*

$$\begin{aligned} \varphi &: \left(\sum_{a:A} P(a) \right) \rightarrow \left(\sum_{b:B} Q(b) \right) \\ \varphi(a, u) &:= (f(a), g_a(u)). \end{aligned}$$

Conversely, if φ and each g_a are n -connected, and moreover Q is fiberwise merely inhabited (i.e. we have $\|Q(b)\|$ for all $b : B$), then f is n -connected.

Proof. For any $b : B$ and $v : Q(b)$ we have

$$\begin{aligned} \|\text{fib}_\varphi((b, v))\|_n &\simeq \left\| \sum_{(a:A)} \sum_{(u:P(a))} \sum_{(p:f(a)=b)} p_*(g_a(u)) = v \right\|_n \\ &\simeq \left\| \sum_{(w:\text{fib}_f(b))} \sum_{(u:P(\text{pr}_1(w)))} g_{\text{pr}_1 w}(u) = \text{pr}_2(w)^{-1}{}_*(v) \right\|_n \\ &\simeq \left\| \sum_{w:\text{fib}_f(b)} \text{fib}_{g(\text{pr}_1 w)}(\text{pr}_2(w)^{-1}{}_*(v)) \right\|_n \\ &\simeq \left\| \sum_{w:\text{fib}_f(b)} \left\| \text{fib}_{g(\text{pr}_1 w)}(\text{pr}_2(w)^{-1}{}_*(v)) \right\|_n \right\|_n \\ &\simeq \|\text{fib}_f(b)\|_n \end{aligned}$$

where the transportations along $f(p)$ and $f(p)^{-1}$ are with respect to Q . Therefore, if either is contractible, so is the other.

In particular, if f is n -connected, then $\|\text{fib}_f(b)\|_n$ is contractible for all $b : B$, and hence so is $\|\text{fib}_\varphi((b, v))\|_n$ for all $(b, v) : \sum_{(b:B)} Q(b)$. On the other hand, if φ is n -connected, then $\|\text{fib}_\varphi((b, v))\|_n$ is contractible for all (b, v) , hence so is $\|\text{fib}_f(b)\|_n$ for any $b : B$ such that there exists some $v : Q(b)$. Finally, since contractibility is a mere proposition, it suffices to merely have such a v . \square

The converse direction of Lemma 7.5.12 can fail if Q is not fiberwise merely inhabited. For example, if P and Q are both constant at $\mathbf{0}$, then φ and each g_a are equivalences, but f could be arbitrary.

In the other direction, we have

Lemma 7.5.13. *Let $P, Q : A \rightarrow \mathcal{U}$ be type families and consider a fiberwise transformation*

$$f : \prod_{a:A} (P(a) \rightarrow Q(a))$$

from P to Q . Then the induced map $\text{total}(f) : \sum_{(a:A)} P(a) \rightarrow \sum_{(a:A)} Q(a)$ is n -connected if and only if each $f(a)$ is n -connected.

Of course, the “only if” direction is also a special case of Lemma 7.5.12.

Proof. By Theorem 4.7.6, we have $\text{fib}_{\text{total}(f)}((x, v)) \simeq \text{fib}_{f(x)}(v)$ for each $x : A$ and $v : Q(x)$. Hence $\|\text{fib}_{\text{total}(f)}((x, v))\|_n$ is contractible if and only if $\|\text{fib}_{f(x)}(v)\|_n$ is contractible. \square

Another useful fact about connected maps is that they induce an equivalence on n -truncations:

Lemma 7.5.14. *If $f : A \rightarrow B$ is n -connected, then it induces an equivalence $\|A\|_n \simeq \|B\|_n$.*

Proof. Let c be the proof that f is n -connected. From left to right, we use the map $\|f\|_n : \|A\|_n \rightarrow \|B\|_n$. To define the map from right to left, by the universal property of truncations, it suffices to give a map $\text{back} : B \rightarrow \|A\|_n$. We can define this map as follows:

$$\text{back}(y) := \|\text{pr}_1\|_n(\text{pr}_1(c(y))).$$

By definition, $c(y)$ has type $\text{isContr}(\|\text{fib}_f(y)\|_n)$, so its first component has type $\|\text{fib}_f(y)\|_n$, and we can obtain an element of $\|A\|_n$ from this by projection.

Next, we show that the composites are the identity. In both directions, because the goal is a path in an n -truncated type, it suffices to cover the case of the constructor $| - |_n$.

In one direction, we must show that for all $x : A$,

$$\|\text{pr}_1\|_n(\text{pr}_1(c(f(x)))) = |x|_n.$$

But $\left|(x, \text{refl}_{f(x)})\right|_n : \|\text{fib}_f(f(x))\|_n$, and $c(f(x))$ says that this type is contractible, so

$$\text{pr}_1(c(f(x))) = |(x, \text{refl})|_n.$$

Applying $\|\text{pr}_1\|_n$ to both sides of this equation gives the result.

In the other direction, we must show that for all $y : B$,

$$\|f\|_n(\|\text{pr}_1\|_n(\text{pr}_1(c(y)))) = |y|_n.$$

$\text{pr}_1(c(y))$ has type $\|\text{fib}_f(y)\|_n$, and the path we want is essentially the second component of the $\text{fib}_f(y)$, but we need to make sure the truncations work out.

In general, suppose we are given $p : \left\|\sum_{(x:A)} B(x)\right\|_n$ and wish to prove $P(\|\text{pr}_1\|_n(p))$. By truncation induction, it suffices to prove $P(|a|_n)$ for all $a : A$ and $b : B(a)$. Applying this principle in this case, it suffices to prove

$$\|f\|_n(|a|_n) = |y|_n$$

given $a : A$ and $b : f(a) = y$. But the left-hand side equals $|f(a)|_n$, so applying $| - |_n$ to both sides of b gives the result. \square

One might guess that this fact characterizes the n -connected maps, but in fact being n -connected is a bit stronger than this. For instance, the inclusion $0_2 : \mathbf{1} \rightarrow \mathbf{2}$ induces an equivalence on (-1) -truncations, but is not surjective (i.e. (-1) -connected). In ?? we will see that the difference in general is an analogous extra bit of surjectivity.

7.6 Orthogonal factorization

In set theory, the surjections and the injections form a unique factorization system: every function factors essentially uniquely as a surjection followed by an injection. We have seen that surjections generalize naturally to n -connected maps, so it is natural to inquire whether these also participate in a factorization system. Here is the corresponding generalization of injections.

Definition 7.6.1. A function $f : A \rightarrow B$ is **n -truncated** if the fiber $\text{fib}_f(b)$ is an n -type for all $b : B$.

In particular, f is (-2) -truncated if and only if it is an equivalence. And of course, A is an n -type if and only if $A \rightarrow \mathbf{1}$ is n -truncated. Moreover, n -truncated maps could equivalently be defined recursively, like n -types.

Lemma 7.6.2. For any $n \geq -2$, a function $f : A \rightarrow B$ is $(n+1)$ -truncated if and only if for all $x, y : A$, the map $\text{ap}_f : (x = y) \rightarrow (f(x) = f(y))$ is n -truncated. In particular, f is (-1) -truncated if and only if it is an embedding in the sense of §4.6.

Proof. Note that for any $(x, p), (y, q) : \text{fib}_f(b)$, we have

$$\begin{aligned} ((x, p) = (y, q)) &= \sum_{r:x=y} (p = \text{ap}_f(r) \bullet q) \\ &= \sum_{r:x=y} (\text{ap}_f(r) = p \bullet q^{-1}) \\ &= \text{fib}_{\text{ap}_f}(p \bullet q^{-1}). \end{aligned}$$

Thus, any path space in any fiber of f is a fiber of ap_f . On the other hand, choosing $b := f(y)$ and $q := \text{refl}_{f(y)}$ we see that any fiber of ap_f is a path space in a fiber of f . The result follows, since f is $(n+1)$ -truncated if all path spaces of its fibers are n -types. \square

We can now construct the factorization, in a fairly obvious way.

Definition 7.6.3. Let $f : A \rightarrow B$ be a function. The **n -image** of f is defined as

$$\text{im}_n(f) := \sum_{b:B} \|\text{fib}_f(b)\|_n.$$

When $n = -1$, we write simply $\text{im}(f)$ and call it the **image** of f .

Lemma 7.6.4. For any function $f : A \rightarrow B$, the canonical function $\tilde{f} : A \rightarrow \text{im}_n(f)$ is n -connected. Consequently, any function factors as an n -connected function followed by an n -truncated function.

Proof. Note that $A \simeq \sum_{(b:B)} \text{fib}_f(b)$. The function \tilde{f} is the function on total spaces induced by the canonical fiberwise transformation

$$\prod_{b:B} (\text{fib}_f(b) \rightarrow \|\text{fib}_f(b)\|_n).$$

Since each map $\text{fib}_f(b) \rightarrow \|\text{fib}_f(b)\|_n$ is n -connected by Corollary 7.5.8, \tilde{f} is n -connected by Lemma 7.5.13. Finally, the projection $\text{pr}_1 : \text{im}_n(f) \rightarrow B$ is n -truncated, since its fibers are equivalent to the n -truncations of the fibers of f . \square

In the following lemma we set up some machinery to prove the unique factorization theorem.

Lemma 7.6.5. Suppose we have a commutative diagram of functions

$$\begin{array}{ccc} A & \xrightarrow{g_1} & X_1 \\ g_2 \downarrow & & \downarrow h_1 \\ X_2 & \xrightarrow{h_2} & B \end{array}$$

with $H : h_1 \circ g_1 \sim h_2 \circ g_2$, where g_1 and g_2 are n -connected and where h_1 and h_2 are n -truncated. Then there is an equivalence

$$E(H, b) : \text{fib}_{h_1}(b) \simeq \text{fib}_{h_2}(b)$$

for any $b : B$, such that for any $a : A$ we have an identification

$$\overline{E}(H, a) : E(H, h_1(g_1(a)))(g_1(a), \text{refl}_{h_1(g_1(a))}) = (g_2(a), H(a)^{-1}).$$

Proof. Let $b : B$. Then we have the following equivalences:

$$\begin{aligned} \text{fib}_{h_1}(b) &\simeq \sum_{w:\text{fib}_{h_1}(b)} \|\text{fib}_{g_1}(\text{pr}_1 w)\|_n && (\text{since } g_1 \text{ is } n\text{-connected}) \\ &\simeq \left\| \sum_{w:\text{fib}_{h_1}(b)} \text{fib}_{g_1}(\text{pr}_1 w) \right\|_n && (\text{by Corollary 7.3.10, since } h_1 \text{ is } n\text{-truncated}) \\ &\simeq \|\text{fib}_{h_1 \circ g_1}(b)\|_n && (\text{by Exercise 4.4}) \end{aligned}$$

and likewise for h_2 and g_2 . Also, since we have a homotopy $H : h_1 \circ g_1 \sim h_2 \circ g_2$, there is an obvious equivalence $\text{fib}_{h_1 \circ g_1}(b) \simeq \text{fib}_{h_2 \circ g_2}(b)$. Hence we obtain

$$\text{fib}_{h_1}(b) \simeq \text{fib}_{h_2}(b)$$

for any $b : B$. By analyzing the underlying functions, we get the following representation of what happens to the element $(g_1(a), \text{refl}_{h_1(g_1(a))})$ after applying each of the equivalences of which E is composed. Some of the identifications are definitional, but others (marked with a = below) are only propositional; putting them together we obtain $\bar{E}(H, a)$.

$$\begin{aligned} (g_1(a), \text{refl}_{h_1(g_1(a))}) &\stackrel{=}{\mapsto} \left((g_1(a), \text{refl}_{h_1(g_1(a))}), \left| (a, \text{refl}_{g_1(a)}) \right|_n \right) \\ &\mapsto \left| ((g_1(a), \text{refl}_{h_1(g_1(a))}), (a, \text{refl}_{g_1(a)})) \right|_n \\ &\mapsto \left| (a, \text{refl}_{h_1(g_1(a))}) \right|_n \\ &\stackrel{=}{\mapsto} \left| (a, H(a)^{-1}) \right|_n \\ &\mapsto \left| ((g_2(a), H(a)^{-1}), (a, \text{refl}_{g_2(a)})) \right|_n \\ &\mapsto \left((g_2(a), H(a)^{-1}), \left| (a, \text{refl}_{g_2(a)}) \right|_n \right) \\ &\mapsto (g_2(a), H(a)^{-1}) \end{aligned}$$

The first equality is because for general b , the map $\text{fib}_{h_1}(b) \rightarrow \sum_{w:\text{fib}_{h_1}(b)} \|\text{fib}_{g_1}(\text{pr}_1 w)\|_n$ inserts the center of contraction for $\|\text{fib}_{g_1}(\text{pr}_1 w)\|_n$ supplied by the assumption that g_1 is n -truncated; whereas in the case in question this type has the obvious inhabitant $\left| (a, \text{refl}_{g_1(a)}) \right|_n$, which by contractibility must be equal to the center. The second propositional equality is because the equivalence $\text{fib}_{h_1 \circ g_1}(b) \simeq \text{fib}_{h_2 \circ g_2}(b)$ concatenates the second components with $H(a)^{-1}$, and we have $H(a)^{-1} \cdot \text{refl} = H(a)^{-1}$. The reader may check that the other equalities are definitional (assuming a reasonable solution to Exercise 4.4). \square

Combining Lemmas 7.6.4 and 7.6.5, we have the following unique factorization result:

Theorem 7.6.6. *For each $f : A \rightarrow B$, the space $\text{fact}_n(f)$ defined by*

$$\sum_{(X:\mathcal{U})} \sum_{(g:A \rightarrow X)} \sum_{(h:X \rightarrow B)} (h \circ g \sim f) \times \text{conn}_n(g) \times \text{trunc}_n(h)$$

is contractible. Its center of contraction is the element

$$(\text{im}_n(f), \tilde{f}, \text{pr}_1, \theta, \varphi, \psi) : \text{fact}_n(f)$$

arising from Lemma 7.6.4, where $\theta : \text{pr}_1 \circ \tilde{f} \sim f$ is the canonical homotopy, where φ is the proof of Lemma 7.6.4, and where ψ is the obvious proof that $\text{pr}_1 : \text{im}_n(f) \rightarrow B$ has n -truncated fibers.

Proof. By Lemma 7.6.4 we know that there is an element of $\text{fact}_n(f)$, hence it is enough to show that $\text{fact}_n(f)$ is a mere proposition. Suppose we have two n -factorizations

$$(X_1, g_1, h_1, H_1, \varphi_1, \psi_1) \quad \text{and} \quad (X_2, g_2, h_2, H_2, \varphi_2, \psi_2)$$

of f . Then we have the pointwise-concatenated homotopy

$$H := (\lambda a. H_1(a) \bullet H_2^{-1}(a)) : (h_1 \circ g_1 \sim h_2 \circ g_2).$$

By univalence and the characterization of paths and transport in Σ -types, function types, and path types, it suffices to show that

- (i) there is an equivalence $e : X_1 \simeq X_2$,
- (ii) there is a homotopy $\zeta : e \circ g_1 \sim g_2$,
- (iii) there is a homotopy $\eta : h_2 \circ e \sim h_1$,
- (iv) for any $a : A$ we have $\text{ap}_{h_2}(\zeta(a))^{-1} \bullet \eta(g_1(a)) \bullet H_1(a) = H_2(a)$.

We prove these four assertions in that order.

- (i) By Lemma 7.6.5, we have a fiberwise equivalence

$$E(H) : \prod_{b:B} \text{fib}_{h_1}(b) \simeq \text{fib}_{h_2}(b).$$

This induces an equivalence of total spaces, i.e. we have

$$\left(\sum_{b:B} \text{fib}_{h_1}(b) \right) \simeq \left(\sum_{b:B} \text{fib}_{h_2}(b) \right).$$

Of course, we also have the equivalences $X_1 \simeq \sum_{(b:B)} \text{fib}_{h_1}(b)$ and $X_2 \simeq \sum_{(b:B)} \text{fib}_{h_2}(b)$ from Lemma 4.8.2. This gives us our equivalence $e : X_1 \simeq X_2$; the reader may verify that the underlying function of e is given by

$$e(x) \equiv \text{pr}_1(E(H, h_1(x))(x, \text{refl}_{h_1(x)})).$$

- (ii) By Lemma 7.6.5, we may choose $\zeta(a) := \text{ap}_{\text{pr}_1}(\bar{E}(H, a)) : e(g_1(a)) = g_2(a)$.
- (iii) For every $x : X_1$, we have

$$\text{pr}_2(E(H, h_1(x))(x, \text{refl}_{h_1(x)})) : h_2(e(x)) = h_1(x),$$

giving us a homotopy $\eta : h_2 \circ e \sim h_1$.

- (iv) By the characterization of paths in fibers (Lemma 4.2.5), the path $\bar{E}(H, a)$ from Lemma 7.6.5 gives us $\eta(g_1(a)) = \text{ap}_{h_2}(\zeta(a)) \bullet H(a)^{-1}$. The desired equality follows by substituting the definition of H and rearranging paths. \square

By standard arguments, this yields the following orthogonality principle.

Theorem 7.6.7. *Let $e : A \rightarrow B$ be n -connected and $m : C \rightarrow D$ be n -truncated. Then the map*

$$\varphi : (B \rightarrow C) \rightarrow \sum_{(h:A \rightarrow C)} \sum_{(k:B \rightarrow D)} (m \circ h \sim k \circ e)$$

is an equivalence.

Sketch of proof. For any (h, k, H) in the codomain, let $h = h_2 \circ h_1$ and $k = k_2 \circ k_1$, where h_1 and k_1 are n -connected and h_2 and k_2 are n -truncated. Then $f = (m \circ h_2) \circ h_1$ and $f = k_2 \circ (k_1 \circ e)$ are both n -factorizations of $m \circ h = k \circ e$. Thus, there is a unique equivalence between them. It is straightforward (if a bit tedious) to extract from this that $\text{fib}_\varphi((h, k, H))$ is contractible. \square

We end by showing that images are stable under pullback.

Lemma 7.6.8. *Suppose that the square*

$$\begin{array}{ccc} A & \longrightarrow & C \\ f \downarrow & & \downarrow g \\ B & \xrightarrow{h} & D \end{array}$$

is a pullback square and let $b : B$. Then $\text{fib}_f(b) \simeq \text{fib}_g(h(b))$.

Proof. This follows from pasting of pullbacks (Exercise 2.12), since the type X in the diagram

$$\begin{array}{ccccc} X & \longrightarrow & A & \longrightarrow & C \\ \downarrow & & f \downarrow & & \downarrow g \\ \mathbf{1} & \xrightarrow{b} & B & \xrightarrow{h} & D \end{array}$$

is the pullback of the left square if and only if it is the pullback of the outer rectangle, while $\text{fib}_f(b)$ is the pullback of the square on the left and $\text{fib}_g(h(b))$ is the pullback of the outer rectangle. \square

Theorem 7.6.9. *Consider functions $f : A \rightarrow B$, $g : C \rightarrow D$ and the diagram*

$$\begin{array}{ccccc} A & \longrightarrow & C & & \\ \tilde{f}_n \downarrow & & \downarrow \tilde{g}_n & & \\ \text{im}_n(f) & \longrightarrow & \text{im}_n(g) & & \\ \text{pr}_1 \downarrow & & \downarrow \text{pr}_1 & & \\ B & \xrightarrow{h} & D & & \end{array}$$

If the outer rectangle is a pullback, then so is the bottom square (and hence so is the top square, by Exercise 2.12). Consequently, images are stable under pullbacks.

Proof. Assuming the outer square is a pullback, we have equivalences

$$\begin{aligned} B \times_D \text{im}_n(g) &\equiv \sum_{(b:B)} \sum_{(w:\text{im}_n(g))} h(b) = \text{pr}_1 w \\ &\simeq \sum_{(b:B)} \sum_{(d:D)} \sum_{(w:\|\text{fib}_g(d)\|_n)} h(b) = d \\ &\simeq \sum_{b:B} \|\text{fib}_g(h(b))\|_n \\ &\simeq \sum_{b:B} \|\text{fib}_f(b)\|_n && \text{(by Lemma 7.6.8)} \\ &\equiv \text{im}_n(f). \end{aligned}$$

\square

7.7 Modalities

Nearly all of the theory of n -types and connectedness can be done in much greater generality. This section will not be used in the rest of the book.

Our first thought regarding generalizing the theory of n -types might be to take Lemma 7.3.3 as a definition.

Definition 7.7.1. A **reflective subuniverse** is a predicate $P : \mathcal{U} \rightarrow \text{Prop}$ such that for every $A : \mathcal{U}$ we have a type $\circlearrowleft A$ such that $P(\circlearrowleft A)$ and a map $\eta_A : A \rightarrow \circlearrowleft A$, with the property that for every $B : \mathcal{U}$ with $P(B)$, the following map is an equivalence:

$$\begin{cases} (\circlearrowleft A \rightarrow B) & \longrightarrow (A \rightarrow B) \\ f & \longmapsto f \circ \eta_A \end{cases} .$$

We write $\mathcal{U}_P := \{ A : \mathcal{U} \mid P(A) \}$, so $A : \mathcal{U}_P$ means that $A : \mathcal{U}$ and we have $P(A)$. We also write $\text{rec}_\circlearrowleft$ for the quasi-inverse of the above map. The notation \circlearrowleft may seem slightly odd, but it will make more sense soon.

For any reflective subuniverse, we can prove all the familiar facts about reflective subcategories from category theory, in the usual way. For instance, we have:

- A type A lies in \mathcal{U}_P if and only if $\eta_A : A \rightarrow \circlearrowleft A$ is an equivalence.
- \mathcal{U}_P is closed under retracts. In particular, A lies in \mathcal{U}_P as soon as η_A admits a retraction.
- The operation \circlearrowleft is a functor in a suitable up-to-coherent-homotopy sense, which we can make precise at as high levels as necessary.
- The types in \mathcal{U}_P are closed under all limits such as products and pullbacks. In particular, for any $A : \mathcal{U}_P$ and $x, y : A$, the identity type $(x =_A y)$ is also in \mathcal{U}_P , since it is a pullback of two functions $\mathbf{1} \rightarrow A$.
- Colimits in \mathcal{U}_P can be constructed by applying \circlearrowleft to ordinary colimits of types.

Importantly, closure under products extends also to “infinite products”, i.e. dependent function types.

Theorem 7.7.2. If $B : A \rightarrow \mathcal{U}_P$ is any family of types in a reflective subuniverse \mathcal{U}_P , then $\prod_{(x:A)} B(x)$ is also in \mathcal{U}_P .

Proof. For any $x : A$, consider the function $\text{ev}_x : (\prod_{(x:A)} B(x)) \rightarrow B(x)$ defined by $\text{ev}_x(f) := f(x)$. Since $B(x)$ lies in P , this extends to a function

$$\text{rec}_\circlearrowleft(\text{ev}_x) : \circlearrowleft\left(\prod_{x:A} B(x)\right) \rightarrow B(x).$$

Thus we can define $h : \circlearrowleft(\prod_{(x:A)} B(x)) \rightarrow \prod_{(x:A)} B(x)$ by $h(z)(x) := \text{rec}_\circlearrowleft(\text{ev}_x)(z)$. Then h is a retraction of $\eta_{\prod_{(x:A)} B(x)}$, so that $\prod_{(x:A)} B(x)$ is in \mathcal{U}_P . \square

In particular, if $B : \mathcal{U}_P$ and A is any type, then $(A \rightarrow B)$ is in \mathcal{U}_P . In categorical language, this means that any reflective subuniverse is an **exponential ideal**. This, in turn, implies by a standard argument that the reflector preserves finite products.

Corollary 7.7.3. For any types A and B and any reflective subuniverse, the induced map $\circlearrowleft(A \times B) \rightarrow \circlearrowleft(A) \times \circlearrowleft(B)$ is an equivalence.

Proof. It suffices to show that $\circ(A) \times \circ(B)$ has the same universal property as $\circ(A \times B)$. It lies in \mathcal{U}_P by the above remark that types in \mathcal{U}_P are closed under limits. Now let $C : \mathcal{U}_P$; we have

$$\begin{aligned} (\circ(A) \times \circ(B) \rightarrow C) &= (\circ(A) \rightarrow (\circ(B) \rightarrow C)) \\ &= (\circ(A) \rightarrow (B \rightarrow C)) \\ &= (A \rightarrow (B \rightarrow C)) \\ &= (A \times B \rightarrow C) \end{aligned}$$

using the universal properties of $\circ(B)$ and $\circ(A)$, along with the fact that $B \rightarrow C$ is in \mathcal{U}_P since C is. It is straightforward to verify that this equivalence is given by composing with $\eta_A \times \eta_B$, as needed. \square

It may seem odd that every reflective subcategory of types is automatically an exponential ideal, with a product-preserving reflector. However, this is also the case classically in the category of *sets*, for the same reasons. It's just that this fact is not usually remarked on, since the classical category of sets—in contrast to the category of homotopy types—does not have many interesting reflective subcategories.

Two basic properties of n -types are *not* shared by general reflective subuniverses: Theorem 7.1.8 (closure under Σ -types) and Theorem 7.3.2 (truncation induction). However, the analogues of these two properties are equivalent to each other.

Theorem 7.7.4. *For a reflective subuniverse \mathcal{U}_P , the following are logically equivalent.*

- (i) *If $A : \mathcal{U}_P$ and $B : A \rightarrow \mathcal{U}_P$, then $\sum_{(x:A)} B(x)$ is in \mathcal{U}_P .*
- (ii) *for every $A : \mathcal{U}$, type family $B : \circ A \rightarrow \mathcal{U}_P$, and map $g : \prod_{(a:A)} B(\eta(a))$, there exists $f : \prod_{(z:\circ A)} B(z)$ such that $f(\eta(a)) = g(a)$ for all $a : A$.*

Proof. Suppose (i). Then in the situation of (ii), the type $\sum_{(z:\circ A)} B(z)$ lies in \mathcal{U}_P , and we have $g' : A \rightarrow \sum_{(z:\circ A)} B(z)$ defined by $g'(a) := (\eta(a), g(a))$. Thus, we have $\text{rec}_{\circ}(g') : \circ A \rightarrow \sum_{(z:\circ A)} B(z)$ such that $\text{rec}_{\circ}(g')(\eta(a)) = (\eta(a), g(a))$.

Now consider the functions $\text{pr}_1 \circ \text{rec}_{\circ}(g') : \circ A \rightarrow \circ A$ and $\text{id}_{\circ A}$. By assumption, these become equal when precomposed with η . Thus, by the universal property of \circ , they are equal already, i.e. we have $p_z : \text{pr}_1(\text{rec}_{\circ}(g')(z)) = z$ for all z . Now we can define $f(z) := p_{z_*}(\text{pr}_2(\text{rec}_{\circ}(g')(z)))$. Using the adjunction property of the equivalence of definition 7.7.1, one can show that the first component of $\text{rec}_{\circ}(g')(\eta(a)) = (\eta(a), g(a))$ is equal to $p_{\eta(a)}$. Thus, its second component yields $f(\eta(a)) = g(a)$, as needed.

Conversely, suppose (ii), and that $A : \mathcal{U}_P$ and $B : A \rightarrow \mathcal{U}_P$. Let h be the composite

$$\circ\left(\sum_{x:A} B(x)\right) \xrightarrow{\circ(\text{pr}_1)} \circ A \xrightarrow{(\eta_A)^{-1}} A.$$

Then for $z : \sum_{(x:A)} B(x)$ we have

$$\begin{aligned} h(\eta(z)) &= \eta^{-1}(\circ(\text{pr}_1)(\eta(z))) \\ &= \eta^{-1}(\eta(\text{pr}_1(z))) \\ &= \text{pr}_1(z). \end{aligned}$$

Denote this path by p_z . Now if we define $C : \circ(\sum_{(x:A)} B(x)) \rightarrow \mathcal{U}$ by $C(w) := B(h(w))$, we have

$$g := \lambda z. p_{z_*}(\text{pr}_2(z)) : \prod_{z:\sum_{(x:A)} B(x)} C(\eta(z)).$$

Thus, the assumption yields $f : \prod_{(w:\bigcirc(\sum_{(x:A)} B(x)))} C(w)$ such that $f(\eta(z)) = g(z)$. Together, h and f give a function $k : \bigcirc(\sum_{(x:A)} B(x)) \rightarrow \sum_{(x:A)} B(x)$ defined by $k(w) := (h(w), f(w))$, while p_z and the equality $f(\eta(z)) = g(z)$ show that k is a retraction of $\eta_{\sum_{(x:A)} B(x)}$. Therefore, $\sum_{(x:A)} B(x)$ is in \mathcal{U}_P . \square

Note the similarity to the discussion in §5.5. The universal property of the reflector of a reflective subuniverse is like a recursion principle with its uniqueness property, while Theorem 7.7.4(ii) is like the corresponding induction principle. Unlike in §5.5, the two are not equivalent here, because of the restriction that we can only eliminate into types that lie in \mathcal{U}_P . Condition (i) of Theorem 7.7.4 is what fixes the disconnect.

Unsurprisingly, of course, if we have the induction principle, then we can derive the recursion principle. We can also derive its uniqueness property, as long as we allow ourselves to eliminate into path types. This suggests the following definition. Note that any reflective subuniverse can be characterized by the operation $\bigcirc : \mathcal{U} \rightarrow \mathcal{U}$ and the functions $\eta_A : A \rightarrow \bigcirc A$, since we have $P(A) = \text{isequiv}(\eta_A)$.

Definition 7.7.5. A **modality** is an operation $\bigcirc : \mathcal{U} \rightarrow \mathcal{U}$ for which there are

- (i) functions $\eta_A^\bigcirc : A \rightarrow \bigcirc(A)$ for every type A .
- (ii) for every $A : \mathcal{U}$ and every type family $B : \bigcirc(A) \rightarrow \mathcal{U}$, a function

$$\text{ind}_\bigcirc : \left(\prod_{a:A} \bigcirc(B(\eta_A^\bigcirc(a))) \right) \rightarrow \prod_{z:\bigcirc(A)} \bigcirc(B(z)).$$

- (iii) A path $\text{ind}_\bigcirc(f)(\eta_A^\bigcirc(a)) = f(a)$ for each $f : \prod_{(a:A)} \bigcirc(B(\eta_A^\bigcirc(a)))$.
- (iv) For any $z, z' : \bigcirc(A)$, the function $\eta_{z=z'}^\bigcirc : (z = z') \rightarrow \bigcirc(z = z')$ is an equivalence.

We say that A is **modal** for \bigcirc if $\eta_A^\bigcirc : A \rightarrow \bigcirc(A)$ is an equivalence, and we write

$$\mathcal{U}_\bigcirc := \{ X : \mathcal{U} \mid X \text{ is } \bigcirc\text{-modal} \} \quad (7.7.6)$$

for the type of modal types.

Conditions (ii) and (iii) are very similar to Theorem 7.7.4(ii), but phrased using $\bigcirc B(z)$ rather than assuming B to be valued in \mathcal{U}_P . This allows us to state the condition purely in terms of the operation \bigcirc , rather than requiring the predicate $P : \mathcal{U} \rightarrow \text{Prop}$ to be given in advance. (It is not entirely satisfactory, since we still have to refer to P not-so-subtly in clause (iv). We do not know whether (iv) follows from (i)–(iii).) However, the stronger-looking property of Theorem 7.7.4(ii) follows from Definition 7.7.5(ii) and (iii), since for any $C : \bigcirc A \rightarrow \mathcal{U}_\bigcirc$ we have $C(z) \simeq \bigcirc C(z)$, and we can pass back across this equivalence.

As with other induction principles, this implies a universal property.

Theorem 7.7.7. Let A be a type and let $B : \bigcirc(A) \rightarrow \mathcal{U}_\bigcirc$. Then the function

$$(- \circ \eta_A^\bigcirc) : \left(\prod_{z:\bigcirc(A)} B(z) \right) \rightarrow \left(\prod_{a:A} B(\eta_A^\bigcirc(a)) \right)$$

is an equivalence.

Proof. By definition, the operation ind_\circ is a right inverse to $(-\circ\eta_A^\circ)$. Thus, we only need to find a homotopy

$$\prod_{z:\circ(A)} s(z) = \text{ind}_\circ(s \circ \eta_A^\circ)(z)$$

for each $s : \prod_{(z:\circ(A))} B(z)$, exhibiting it as a left inverse as well. By assumption, each $B(z)$ is modal, and hence each type $s(z) = R_X^\circ(s \circ \eta_A^\circ)(z)$ is also modal. Thus, it suffices to find a function of type

$$\prod_{a:A} s(\eta_A^\circ(a)) = \text{ind}_\circ(s \circ \eta_A^\circ)(\eta_A^\circ(a))$$

which follows from Definition 7.7.5(iii). \square

In particular, for every type A and every modal type B , we have an equivalence $(\circ A \rightarrow B) \simeq (A \rightarrow B)$.

Corollary 7.7.8. *For any modality \circ , the \circ -modal types form a reflective subuniverse satisfying the equivalent conditions of Theorem 7.7.4.*

Thus, modalities can be identified with reflective subuniverses closed under Σ -types. The name *modality* comes, of course, from *modal logic*, which studies logic where we can form statements such as “possibly A ” (usually written $\diamond A$) or “necessarily A ” (usually written $\square A$). The symbol \circ is somewhat common for an arbitrary modal operator. Under the propositions-as-types principle, a modality in the sense of modal logic corresponds to an operation on *types*, and Definition 7.7.5 seems a reasonable candidate for how such an operation should be defined. (More precisely, we should perhaps call these *idempotent, monadic* modalities; see the Notes.) As mentioned in §3.10, we may in general use adverbs to speak informally about such modalities, such as “merely” for the propositional truncation and “purely” for the identity modality (i.e. the one defined by $\circ A : \equiv A$).

For any modality \circ , we define a map $f : A \rightarrow B$ to be **\circ -connected** if $\circ(\text{fib}_f(b))$ is contractible for all $b : B$, and to be **\circ -truncated** if $\text{fib}_f(b)$ is modal for all $b : B$. All of the theory of §§7.5 and 7.6 which doesn’t involve relating n -types for different values of n applies verbatim in this generality. In particular, we have an orthogonal factorization system.

An important class of modalities which does *not* include the n -truncations is the *left exact* modalities: those for which the functor \circ preserves pullbacks as well as finite products. These are a categorification of “Lawvere-Tierney topologies” in elementary topos theory, and correspond in higher-categorical semantics to sub- $(\infty, 1)$ -toposes. However, this is beyond the scope of this book.

Some particular examples of modalities other than n -truncation can be found in the exercises.

Notes

The notion of homotopy n -type in classical homotopy theory is quite old. It was Voevodsky who realized that the notion can be defined recursively in homotopy type theory, starting from contractibility.

The property “Axiom K” was so named by Thomas Streicher, as a property of identity types which comes after J, the latter being the traditional name for the eliminator of identity types. Theorem 7.2.5 is due to Hedberg [Hed98]; [KECA13] contains more information and generalizations.

The notions of n -connected spaces and functions are also classical in homotopy theory, although as mentioned before, our indexing for connectedness of functions is off by one from the classical indexing. The importance of the resulting factorization system has been emphasized by recent work in higher topos theory by Rezk, Lurie, and others. In particular, the results of this chapter should be compared with [Lur09, §6.5.1]. In ??, the theory of n -connected maps will be crucial to our proof of the Freudenthal suspension theorem.

Modal operators in *simple* type theory have been studied extensively; see e.g. [dPGM04]. In the setting of dependent type theory, [AB04] treats the special case of propositional truncation ((-1) -truncation) as a modal operator. The development presented here greatly extends and generalizes this work, while drawing also on ideas from topos theory.

Generally, modal operators come in (at least) two flavors: those such as \diamond (“possibly”) for which $A \Rightarrow \diamond A$, and those such as \square (“necessarily”) for which $\square A \Rightarrow A$. When they are also *idempotent* (i.e. $\diamond A = \diamond\diamond A$ or $\square A = \square\square A$), the former may be identified with reflective subcategories (or equivalently, idempotent monads), and the latter with coreflective subcategories (or idempotent comonads). However, in dependent type theory it is trickier to deal with the comonadic sort, since they are more rarely stable under pullback, and thus cannot be interpreted as operations on the universe \mathcal{U} . Sometimes there are ways around this (see e.g. [SS12]), but for simplicity, here we stick to the monadic sort.

On the computational side, monads (and hence modalities) are used to model computational effects in functional programming [Mog89]. A computation is said to be *pure* if its execution results in no side effects (such as printing a message to the screen, playing music, or sending data over the Internet). There exist “purely functional” programming languages, such as Haskell, in which it is technically only possible to write pure functions: side effects are represented by applying “monads” to output types. For instance, a function of type $\text{Int} \rightarrow \text{Int}$ is pure, while a function of type $\text{Int} \rightarrow \text{IO}(\text{Int})$ may perform input and output along the way to computing its result; the operation IO is a monad. (This is the origin of our use of the adverb “purely” for the identity monad, since it corresponds computationally to pure functions with no side-effects.) The modalities we have considered in this chapter are all idempotent, whereas those used in functional programming rarely are, but the ideas are still closely related.

Exercises

Exercise 7.1.

- (i) Use Theorem 7.2.2 to show that if $\|A\| \rightarrow A$ for every type A , then every type is a set.
- (ii) Show that if every surjective function (purely) splits, i.e. if $\prod_{(b:B)} \|\text{fib}_f(b)\| \rightarrow \prod_{(b:B)} \text{fib}_f(b)$ for every $f : A \rightarrow B$, then every type is a set.

Exercise 7.2. For this exercise, we consider the following general notion of colimit. Define a **graph** Γ to consist of a type Γ_0 and a family $\Gamma_1 : \Gamma_0 \rightarrow \Gamma_0 \rightarrow \mathcal{U}$. A **diagram** (of types) over a graph Γ consists of a family $F : \Gamma_0 \rightarrow \mathcal{U}$ together with for each $x, y : \Gamma_0$, a function $F_{x,y} : \Gamma_1(x, y) \rightarrow F(x) \rightarrow F(y)$. The **colimit** of such a diagram is the higher inductive type $\text{colim}(F)$ generated by

- for each $x : \Gamma_0$, a function $\text{inc}_x : F(x) \rightarrow \text{colim}(F)$, and
- for each $x, y : \Gamma_0$ and $\gamma : \Gamma_1(x, y)$ and $a : F(x)$, a path $\text{inc}_y(F_{x,y}(\gamma, a)) = \text{inc}_x(a)$.

There are more general kinds of colimits as well (see e.g. Exercise 7.16), but this is good enough for many purposes.

- (i) Exhibit a graph Γ such that colimits of Γ -diagrams can be identified with pushouts as defined in §6.8. In other words, each span should induce a diagram over Γ whose colimit is the pushout of the span.
- (ii) Exhibit a graph Γ and a diagram F over Γ such that $F(x) = \mathbf{1}$ for all x , but such that $\text{colim}(F) = S^2$. Note that $\mathbf{1}$ is a (-2) -type, while S^2 is not expected to be an n -type for any finite n . See also Exercise 7.16.

Exercise 7.3. Show that if A is an n -type and $B : A \rightarrow n\text{-Type}$ is a family of n -types, where $n \geq -1$, then the W -type $W_{(a:A)} B(a)$ (see §5.3) is also an n -type.

Exercise 7.4. Use Lemma 7.5.13 to extend Lemma 7.5.11 to any section-retraction pair.

Exercise 7.5. Show that Corollary 7.5.9 also works as a characterization in the other direction: B is an n -type if and only if every map into B from an n -connected type is constant. Ideally, your proof should work for any modality as in §7.7.

Exercise 7.6. Prove that for $n \geq -1$, a type A is n -connected if and only if it is merely inhabited and for all $a, b : A$ the type $a =_A b$ is $(n-1)$ -connected. Thus, since every type is (-2) -connected, n -connectedness of types can be defined inductively using only propositional truncations. (In particular, A is 0-connected if and only if $\|A\|$ and $\prod_{(a,b:A)} \|a = b\|$.)

Exercise 7.7. For $-1 \leq n, m \leq \infty$, let $\text{LEM}_{n,m}$ denote the statement

$$\prod_{A:n\text{-Type}} \|A + \neg A\|_m,$$

where $\infty\text{-Type} := \mathcal{U}$ and $\|X\|_\infty := X$. Show that:

- (i) If $n = -1$ or $m = -1$, then $\text{LEM}_{n,m}$ is equivalent to LEM from §3.4.
- (ii) If $n \geq 0$ and $m \geq 0$, then $\text{LEM}_{n,m}$ is inconsistent with univalence.

Exercise 7.8. For $-1 \leq n, m \leq \infty$, let $\text{AC}_{n,m}$ denote the statement

$$\prod_{(X:\text{Set})} \prod_{(Y:X \rightarrow n\text{-Type})} \left(\prod_{x:X} \|Y(x)\|_m \right) \rightarrow \left\| \prod_{x:X} Y(x) \right\|_m,$$

with conventions as in Exercise 7.7. Thus $\text{AC}_{0,-1}$ is the axiom of choice from §3.8, while $\text{AC}_{\infty,\infty}$ is the identity function. (If we had formulated $\text{AC}_{n,m}$ analogously to (3.8.1) rather than (3.8.3), $\text{AC}_{\infty,\infty}$ would be like Theorem 2.15.7.) It is known that $\text{AC}_{\infty,-1}$ is consistent with univalence, since it holds in Voevodsky's simplicial model.

- (i) Without using univalence, show that $\text{LEM}_{n,\infty}$ implies $\text{AC}_{n,m}$ for all m . (On the other hand, in ?? we will show that $\text{AC} = \text{AC}_{0,-1}$ implies $\text{LEM} = \text{LEM}_{-1,-1}$.)
- (ii) Of course, $\text{AC}_{n,m} \Rightarrow \text{AC}_{k,m}$ if $k \leq n$. Are there any other implications between the principles $\text{AC}_{n,m}$? Is $\text{AC}_{n,m}$ consistent with univalence for any $m \geq 0$ and any n ? (These are open questions.)

Exercise 7.9. Show that $\text{AC}_{n,-1}$ implies that for any n -type A , there merely exists a set B and a surjection $B \rightarrow A$.

Exercise 7.10. Define the **n -connected axiom of choice** to be the statement

If X is a set and $Y : X \rightarrow \mathcal{U}$ is a family of types such that each $Y(x)$ is n -connected, then $\prod_{(x:X)} Y(x)$ is n -connected.

Note that the (-1) -connected axiom of choice is $\text{AC}_{\infty, -1}$ from Exercise 7.8.

- (i) Prove that the (-1) -connected axiom of choice implies the n -connected axiom of choice for all $n \geq -1$.
- (ii) Are there any other implications between the n -connected axioms of choice and the principles $\text{AC}_{n,m}$? (This is an open question.)

Exercise 7.11. Show that the n -truncation modality is not left exact for any $n \geq -1$. That is, exhibit a pullback which it fails to preserve.

Exercise 7.12. Show that $X \mapsto (\neg\neg X)$ is a modality.

Exercise 7.13. Let P be a mere proposition.

- (i) Show that $X \mapsto (P \rightarrow X)$ is a left exact modality. This is called the **open modality** associated to P .
- (ii) Show that $X \mapsto P * X$ is a left exact modality, where $*$ denotes the join (see §6.8). This is called the **closed modality** associated to P .

Exercise 7.14. Let $f : A \rightarrow B$ be a map; a type Z is **f -local** if $(-\circ f) : (B \rightarrow Z) \rightarrow (A \rightarrow Z)$ is an equivalence.

- (i) Prove that the f -local types form a reflective subuniverse. You will want to use a higher inductive type to define the reflector (localization).
- (ii) Prove that if $B = \mathbf{1}$, then this subuniverse is a modality.

Exercise 7.15. Show that in contrast to Remark 6.7.1, we could equivalently define $\|A\|_n$ to be generated by a function $|-\|_n : A \rightarrow \|A\|_n$ together with for each $r : S^{n+1} \rightarrow \|A\|_n$ and each $x : S^{n+1}$, a path $s_r(x) : r(x) = r(\text{base})$.

Exercise 7.16. In this exercise, we consider a slightly fancier notion of colimit than in Exercise 7.2. Define a **graph with composition** Γ to be a graph as in Exercise 7.2 together with for each $x, y, z : \Gamma_0$, a function $\Gamma_1(y, z) \rightarrow \Gamma_1(x, y) \rightarrow \Gamma_1(x, z)$, written as $\delta \mapsto \gamma \mapsto \delta \circ \gamma$. (For instance, any precategory as in ?? is a graph with composition.) A **diagram** F over a graph with composition Γ consists of a diagram over the underlying graph, together with for each $x, y, z : \Gamma_0$ and $\gamma : \Gamma_1(x, y)$ and $\delta : \Gamma_1(y, z)$, a homotopy $\text{cmp}_{x,y,z}(\delta, \gamma) : F_{y,z}(\delta) \circ F_{x,y}(\gamma) \sim F_{x,z}(\delta \circ \gamma)$. The **colimit** of such a diagram is the higher inductive type $\text{colim}(F)$ generated by

- for each $x : \Gamma_0$, a function $\text{inc}_x : F(x) \rightarrow \text{colim}(F)$,
- for each $x, y : \Gamma_0$ and $\gamma : \Gamma_1(x, y)$ and $a : F(x)$, a path $\text{glue}_{x,y}(\gamma, a) : \text{inc}_y(F_{x,y}(\gamma, a)) = \text{inc}_x(a)$, and
- for each $x, y, z : \Gamma_0$ and $\gamma : \Gamma_1(x, y)$ and $\delta : \Gamma_1(y, z)$ and $a : F(x)$, a path

$$\text{inc}_z(\text{cmp}_{x,y,z}(\delta, \gamma, a)) \cdot \text{glue}_{x,z}(\delta \circ \gamma, a) = \text{glue}_{y,z}(\delta, F_{x,y}(\gamma, a)) \cdot \text{glue}_{x,y}(\gamma, a).$$

(This is a “second-order approximation” to a fully homotopy-theoretic notions of diagram and colimit, which ought to involve “coherence paths” of this sort at all higher levels. Defining such things in type theory is an important open problem.)

Exhibit a graph with composition Γ such that Γ_0 is a set and each type $\Gamma_1(x, y)$ is a mere proposition, and a diagram F over Γ such that $F(x) = \mathbf{1}$ for all x , for which $\text{colim}(F) = S^2$.

Exercise 7.17. Comparing Lemmas 7.5.12 and 7.5.13, one might be tempted to conjecture that if $f : A \rightarrow B$ is n -connected and $g : \prod_{(a:A)} P(a) \rightarrow Q(f(a))$ induces an n -connected map $(\sum_{(a:A)} P(a)) \rightarrow (\sum_{(b:B)} Q(b))$, then g is fiberwise n -connected. Give a counterexample to show that this is false. (In fact, when generalized to modalities, this property characterizes the left exact ones; see Exercise 7.13.)

Exercise 7.18. Show that if $f : A \rightarrow B$ is n -connected, then $\|f\|_k : \|A\|_k \rightarrow \|B\|_k$ is also n -connected.

Exercise 7.19. We say a type A is **categorically connected** if for every types B, C the canonical map $e_{A,B,C} : ((A \rightarrow B) + (A \rightarrow C)) \rightarrow (A \rightarrow B + C)$ defined by

$$\begin{aligned} e_{A,B,C}(\text{inl}(g)) &:= \lambda x. \text{inl}(g(x)), \\ e_{A,B,C}(\text{inr}(g)) &:= \lambda x. \text{inr}(g(x)) \end{aligned}$$

is an equivalence.

- (i) Show that any connected type is categorically connected.
- (ii) Show that all categorically connected types are connected if and only if LEM holds. (Hint: consider $A := \Sigma P$ such that $\neg\neg P$ holds.)

Bibliography

- [AB04] Steven Awodey and Andrej Bauer. Propositions as [types]. *Journal of Logic and Computation*, 14(4):447–471, 2004. (Cited on pages 117, 203, and 234.)
- [AG02] Peter Aczel and Nicola Gambino. Collection principles in dependent type theory. In Paul Callaghan, Zhaohui Luo, James McKinna, and Robert Pollack, editors, *Types for Proofs and Programs, International Workshop, TYPES 2000, Durham, UK, December 8-12, 2000, Selected Papers*, volume 2277 of *Lecture Notes in Computer Science*, pages 1–23. Springer, 2002. (Cited on page 117.)
- [AGS12] Steve Awodey, Nicola Gambino, and Kristina Sojakova. Inductive types in homotopy type theory. In *Proceedings of the 2012 27th Annual IEEE/ACM Symposium on Logic in Computer Science*, pages 95–104. IEEE Computer Society, 2012, arXiv:1201.3898. (Cited on page 163.)
- [AKL13] Jeremy Avigad, Krzysztof Kapulkin, and Peter LeFanu Lumsdaine. Homotopy limits in Coq, 2013. arXiv:1304.0680. (Cited on page 97.)
- [Alt99] Thorsten Altenkirch. Extensional equality in intensional type theory. In *14th Annual IEEE Symposium on Logic in Computer Science, Trento, Italy, July 2–5, 1999*, pages 412–420, 1999. (Cited on page 204.)
- [AMS07] Thorsten Altenkirch, Conor McBride, and Wouter Swierstra. Observational equality, now! In Aaron Stump and Hongwei Xi, editors, *Proceedings of the ACM Workshop Programming Languages meets Program Verification, PLPV 2007, Freiburg, Germany, October 5, 2007*, 2007. (Cited on page 204.)
- [Ang13] Carlo Angiuli. The $(\infty, 1)$ -accidentopos model of unintentional type theory. *Sigbovik '13*, April 1 2013. (No citations.)
- [AW09] Steve Awodey and Michael A. Warren. Homotopy theoretic models of identity types. *Mathematical Proceedings of the Cambridge Philosophical Society*, 146:45–55, 2009. (Cited on pages 3 and 96.)
- [Bau13] Andrej Bauer. Five stages of accepting constructive mathematics, 2013. <http://video.ias.edu/members/1213/0318-AndrejBauer>. (No citations.)
- [BCH13] Bruno Barras, Thierry Coquand, and Simon Huber. A generalization of Takeuti-Gandy interpretation. <http://uf-ias-2012.wikispaces.com/file/view/semi.pdf>, 2013. (Cited on page 10.)
- [Bee85] Michael Beeson. *Foundations of Constructive Mathematics*. Springer, 1985. (Cited on page 51.)
- [Bou68] Nicolas Bourbaki. *Theory of Sets*. Hermann, Paris, 1968. (Cited on page 96.)
- [CAB⁺86] Robert L. Constable, Stuart F. Allen, H. M. Bromley, W. R. Cleaveland, J. F. Cremer, Robert W. Harper, Douglas J. Howe, T. B. Knoblock, N. P. Mendler, P. Panangaden, James T. Sasaki, and Scott F. Smith. *Implementing Mathematics with the NuPRL Proof Development System*. Prentice Hall, 1986. (Cited on pages 51, 53, 117, and 203.)

-
- [Car95] Aurelio Carboni. Some free constructions in realizability and proof theory. *Journal of Pure and Applied Algebra*, 103:117–148, 1995. (Cited on page 204.)
- [CDP14] Jesper Cockx, Dominique Devriese, and Frank Piessens. Pattern matching without K. In *Proceedings of the 19th ACM SIGPLAN International Conference on Functional Programming, ICFP 2014, Gothenburg, Sweden, September 1-3, 2014*, 2014. (Cited on page 52.)
- [Chu40] Alonzo Church. A formulation of the simple theory of types. *Journal of Symbolic Logic*, 5:56–68, 1940. (Cited on page 2.)
- [Chu41] Alonzo Church. *The Calculi of Lambda Conversation*. Princeton University Press, 1941. (Cited on page 2.)
- [CM85] Robert L. Constable and N. P. Mendler. Recursive definitions in type theory. In Rohit Parikh, editor, *Logics of Programs, Conference, Brooklyn College, June 17–19, 1985, Proceedings*, volume 193 of *Lecture Notes in Computer Science*, pages 61–78, 1985. (Cited on page 163.)
- [Con85] Robert L. Constable. Constructive mathematics as a programming logic I: Some principles of theory. In *Annals of Mathematics*, volume 24, pages 21–37. Elsevier Science Publishers, B.V. (North-Holland), 1985. Reprinted from *Topics in the Theory of Computation, Selected Papers of the International Conference on Foundations of Computation Theory, FCT '83*. (Cited on page 117.)
- [Coq92a] Thierry Coquand. The paradox of trees in type theory. *BIT Numerical Mathematics*, 32(1):10–14, 1992. (Cited on page 23.)
- [Coq92b] Thierry Coquand. Pattern matching with dependent types. In *Proceedings of the Workshop on Types for Proofs and Programs*, pages 71–83, 1992. (Cited on page 52.)
- [Coq12] Coq Development Team. *The Coq Proof Assistant Reference Manual*. INRIA-Rocquencourt, 2012. (Cited on page 51.)
- [CP90] Thierry Coquand and Christine Paulin. Inductively defined types. In *COLOG-88 (Tallinn, 1988)*, volume 416 of *Lecture Notes in Computer Science*, pages 50–66. Springer, 1990. (Cited on page 163.)
- [dB73] Nicolaas Govert de Bruijn. *AUTOMATH, a language for mathematics*. Les Presses de l’Université de Montréal, Montreal, Quebec, 1973. Séminaire de Mathématiques Supérieures, No. 52 (Été 1971). (Cited on pages 51 and 52.)
- [dPGM04] Valeria de Paiva, Rajeev Goré, and Michael Mendler. Modalities in constructive logics and type theories. *Journal of Logic and Computation*, 14(4):439–446, 2004. (Cited on page 234.)
- [Dyb91] Peter Dybjer. Inductive sets and families in Martin-Löf’s type theory and their set-theoretic semantics. In Gerard Huet and Gordon Plotkin, editors, *Logical Frameworks*, pages 280–30. Cambridge University Press, 1991. (Cited on page 163.)
- [Dyb00] Peter Dybjer. A general formulation of simultaneous inductive-recursive definitions in type theory. *Journal of Symbolic Logic*, 65(2):525–549, 2000. (Cited on page 163.)
- [EucBC] Euclid. *Elements*, Vols. 1–13. Elsevier, 300 BC. (Cited on page 57.)
- [GAA⁺13] Georges Gonthier, Andrea Asperti, Jeremy Avigad, Yves Bertot, Cyril Cohen, François Garillot, Stéphane Le Roux, Assia Mahboubi, Russell O’Connor, Sidi Ould Biha, Ioana Pasca, Laurence Rideau, Alexey Solovyev, Enrico Tassi, and Laurent Théry. A machine-checked proof of the odd order theorem. In *Interactive Theorem Proving*, 2013. (Cited on page 6.)
- [Gar09] Richard Garner. On the strength of dependent products in the type theory of Martin-Löf. *Annals of Pure and Applied Logic*, 160(1):1–12, 2009. (Cited on page 96.)
- [Hed98] Michael Hedberg. A coherence theorem for Martin-Löf’s type theory. *Journal of Functional Programming*, 8(4):413–436, 1998. (Cited on page 233.)

- [Hey66] Arend Heyting. *Intuitionism: an introduction*. Studies in logic and the foundations of mathematics. North-Holland Pub. Co., 1966. (Cited on page 51.)
- [Hof95] Martin Hofmann. *Extensional concepts in intensional type theory*. PhD thesis, University of Edinburgh, 1995. (Cited on page 204.)
- [How80] William A. Howard. The formulae-as-types notion of construction. In J. Roger Seldin, Jonathan P.; Hindley, editor, *To H.B. Curry: Essays on Combinatory Logic, Lambda Calculus and Formalism*, pages 479–490. Academic Press, 1980. original paper manuscript from 1969. (Cited on pages 51, 52, and 96.)
- [HS98] Martin Hofmann and Thomas Streicher. The groupoid interpretation of type theory. In Giovanni Sambin and Jan M. Smith, editors, *Twenty-five years of constructive type theory (Venice, 1995)*, volume 36 of *Oxford Logic Guides*, pages 83–111. Oxford University Press, New York, 1998. (Cited on pages 4 and 96.)
- [Jac99] Bart Jacobs. *Categorical logic and type theory*, volume 141 of *Studies in Logic and the Foundations of Mathematics*. Elsevier, 1999. (Cited on page 117.)
- [Joh02] Peter T. Johnstone. *Sketches of an Elephant: A Topos Theory Compendium: Volumes 1 and 2*. Number 43 in Oxford Logic Guides. Oxford Science Publications, 2002. (Cited on page 117.)
- [KECA13] Nicolai Kraus, Martin Escardó, Thierry Coquand, and Thorsten Altenkirch. Generalizations of Hedberg’s theorem. In Masahito Hasegawa, editor, *11th International Conference, Typed Lambda Calculus and Applications 2013, Eindhoven, The Netherlands, June 26–28, 2013. Proceedings*, volume 7941 of *Lecture Notes in Computer Science*, pages 173–188. Springer Berlin Heidelberg, 2013. (Cited on pages 118 and 233.)
- [KLN04] Fairouz Kamareddine, Twan Laan, and Rob Nederpelt. *A Modern Perspective on Type Theory: From its Origins until Today*. Number 29 in Applied Logic. Kluwer, 2004. (Cited on page 2.)
- [KLV12] Chris Kapulkin, Peter LeFanu Lumsdaine, and Vladimir Voevodsky. The simplicial model of univalent foundations, 2012. arXiv:1211.2851. (Cited on pages 10, 96, and 163.)
- [Kol32] Andrey Kolmogorov. Zur Deutung der intuitionistischen Logik. *Mathematische Zeitschrift*, 35:58–65, 1932. (Cited on page 8.)
- [Law05] F. William Lawvere. An elementary theory of the category of sets (long version) with commentary. *Reprints in Theory and Applications of Categories*, 11:1–35, 2005. Reprinted and expanded from Proc. Nat. Acad. Sci. U.S.A. **52** (1964), With comments by the author and Colin McLarty. (Cited on page 6.)
- [Law06] F. William Lawvere. Adjointness in foundations. *Reprints in Theory and Applications of Categories*, 16:1–16, 2006. Reprinted from *Dialectica* **23** (1969). (Cited on pages 51 and 163.)
- [LH12] Daniel R. Licata and Robert Harper. Canonicity for 2-dimensional type theory. In *Proceedings of the 39th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 337–348, New York, NY, USA, 2012. ACM. (Cited on pages 9, 10, and 96.)
- [LS13] Daniel R. Licata and Michael Shulman. Calculating the fundamental group of the circle in homotopy type theory. In *LICS 2013: Proceedings of the Twenty-Eighth Annual ACM/IEEE Symposium on Logic in Computer Science*, 2013. (Cited on page 96.)
- [LS17] Peter LeFanu Lumsdaine and Michael Shulman. Semantics of higher inductive types. arXiv:1705.07088, 2017. (Cited on pages 10, 163, and 203.)
- [Lum10] Peter LeFanu Lumsdaine. Weak ω -categories from intensional type theory. *Typed lambda calculi and applications*, 6:1–19, 2010. arXiv:0812.0409. (Cited on page 96.)
- [Lur09] Jacob Lurie. *Higher topos theory*. Number 170 in Annals of Mathematics Studies. Princeton University Press, 2009. arXiv:math.CT/0608040. (Cited on pages 9, 137, and 234.)

-
- [ML71] Per Martin-Löf. *Hauptsatz for the intuitionistic theory of iterated inductive definitions*. In *Proceedings of the Second Scandinavian Logic Symposium (University of Oslo 1970)*, volume 63 of *Studies in Logic and the Foundations of Mathematics*, pages 179–216. North-Holland, 1971. (Cited on page 163.)
- [ML75] Per Martin-Löf. An intuitionistic theory of types: predicative part. In H.E. Rose and J.C. Shepherdson, editors, *Logic Colloquium '73, Proceedings of the Logic Colloquium*, volume 80 of *Studies in Logic and the Foundations of Mathematics*, pages 73–118. North-Holland, 1975. (Cited on pages 2, 51, and 163.)
- [ML82] Per Martin-Löf. Constructive mathematics and computer programming. In L. Jonathan Cohen, Jerzy Łoś, Helmut Pfeiffer, and Klaus-Peter Podewski, editors, *Logic, Methodology and Philosophy of Science VI, Proceedings of the Sixth International Congress of Logic, Methodology and Philosophy of Science, Hannover 1979*, volume 104 of *Studies in Logic and the Foundations of Mathematics*, pages 153–175. North-Holland, 1982. (Cited on pages 2, 51, and 163.)
- [ML84] Per Martin-Löf. *Intuitionistic type theory*, volume 1 of *Studies in Proof Theory*. Bibliopolis, 1984. Notes by Giovanni Sambin of a series of lectures given in Padua, June 1980. (Cited on pages 2, 51, and 52.)
- [ML98] Per Martin-Löf. An intuitionistic theory of types. In Giovanni Sambin and Jan M. Smith, editors, *Twenty-five years of constructive type theory (Venice, 1995)*, volume 36 of *Oxford Logic Guides*, pages 127–172. Oxford University Press, 1998. (Cited on pages 2, 51, 52, 53, and 95.)
- [ML06] Per Martin-Löf. 100 years of Zermelo’s axiom of choice: what was the problem with it? *The Computer Journal*, 49(3):345–350, 2006. (Cited on page 117.)
- [Mog89] Eugenio Moggi. Notions of computation and monads. *Information and Computation*, 93:55–92, 1989. (Cited on page 234.)
- [MP00] Ieke Moerdijk and Erik Palmgren. Wellfounded trees in categories. In *Proceedings of the Workshop on Proof Theory and Complexity, PTAC'98 (Aarhus)*, volume 104, pages 189–218, 2000. (Cited on page 163.)
- [MS05] Maria Emilia Maietti and Giovanni Sambin. Toward a minimalist foundation for constructive mathematics. In Laura Crosilla and Peter Schuster, editors, *From Sets and Types to Topology and Analysis: Practicable Foundations for Constructive Mathematics*, volume 48 of *Oxford Logic Guides*, pages 91–114. Clarendon Press, 2005. (Cited on page 117.)
- [Nor07] Ulf Norell. *Towards a practical programming language based on dependent type theory*. PhD thesis, Chalmers, Göteborg University, 2007. (Cited on page 51.)
- [Pie02] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002. (Cited on page 2.)
- [PM93] Christine Paulin-Mohring. Inductive Definitions in the System Coq - Rules and Properties. In Marc Bezem and Jan Friso Groote, editors, *Proceedings of the conference Typed Lambda Calculi and Applications*, number 664 in Lecture Notes in Computer Science, 1993. (Cited on page 53.)
- [PPM90] Frank Pfenning and Christine Paulin-Mohring. Inductively defined types in the calculus of constructions. In Michael G. Main, Austin Melton, Michael W. Mislove, and David A. Schmidt, editors, *Mathematical Foundations of Programming Semantics, 5th International Conference, Tulane University, New Orleans, Louisiana, USA, March 29 – April 1, 1989, Proceedings*, number 442 in Lecture Notes in Computer Science, pages 209–228. Springer, 1990. (Cited on page 163.)
- [PS89] Kent Petersson and Dan Synek. A set constructor for inductive sets in Martin-Löf’s type theory. In David H. Pitt, David E. Rydeheard, Peter Dybjer, Andrew M. Pitts, and Axel Poigné, editors, *Category Theory and Computer Science, Manchester, UK, September 5–8, 1989, Proceedings*, volume 389 of *Lecture Notes in Computer Science*, pages 128–140. Springer, 1989. (Cited on page 163.)

- [Rez05] Charles Rezk. Toposes and homotopy toposes. <http://www.math.uiuc.edu/~rezk/homotopy-topos-sketch.pdf>, 2005. (Cited on pages 9 and 137.)
- [Rus08] Bertrand Russell. Mathematical logic based on the theory of types. *American Journal of Mathematics*, 30:222–262, 1908. (Cited on page 2.)
- [Sco70] Dana Scott. Constructive validity. In M. Laudet, D. Lacombe, L. Nolin, and M. Schützenberger, editors, *Symposium on Automatic Demonstration*, volume 125, pages 237–275. Springer-Verlag, 1970. (Cited on page 51.)
- [Som10] Giovanni Sommaruga. *History and Philosophy of Constructive Type Theory*. Number 290 in Synthese Library. Kluwer, 2010. (Cited on page 2.)
- [Spi11] Arnaud Spiwack. *A Journey Exploring the Power and Limits of Dependent Type Theory*. PhD thesis, École Polytechnique, Palaiseau, France, 2011. (Cited on page 117.)
- [SS12] Urs Schreiber and Michael Shulman. Quantum gauge field theory in cohesive homotopy type theory. *Quantum Physics and Logic*, 2012. (Cited on page 234.)
- [Str93] Thomas Streicher. Investigations into intensional type theory, 1993. Habilitationsschrift, Ludwig-Maximilians-Universität München. (Cited on pages 52, 53, and 208.)
- [Tai67] William W. Tait. Intensional interpretations of functionals of finite type. I. *The Journal of Symbolic Logic*, 32:198–212, 1967. (Cited on page 51.)
- [Tai68] William W. Tait. Constructive reasoning. In *Logic, Methodology and Philos. Sci. III (Proc. Third Internat. Congr., Amsterdam, 1967)*, pages 185–199. North-Holland, Amsterdam, 1968. (Cited on pages 51 and 52.)
- [TV02] Bertrand Toën and Gabriele Vezzosi. Homotopical algebraic geometry I: Topos theory, 2002. arXiv:math/0207028. (Cited on page 9.)
- [Tvd88a] Anne Sjerp Troelstra and Dirk van Dalen. *Constructivism in mathematics. Vol. I*, volume 121 of *Studies in Logic and the Foundations of Mathematics*. North-Holland Publishing Co., Amsterdam, 1988. An introduction. (Cited on pages 8 and 117.)
- [Tvd88b] Anne Sjerp Troelstra and Dirk van Dalen. *Constructivism in mathematics. Vol. II*, volume 123 of *Studies in Logic and the Foundations of Mathematics*. North-Holland Publishing Co., Amsterdam, 1988. An introduction. (Cited on pages 8 and 117.)
- [vdBG11] Benno van den Berg and Richard Garner. Types are weak ω -groupoids. *Proceedings of the London Mathematical Society*, 102(2):370–394, 2011, <http://plms.oxfordjournals.org/content/102/2/370.full.pdf+html>. (Cited on page 96.)
- [vdBM15] Benno van den Berg and Ieke Moerdijk. W-types in homotopy type theory. *Mathematical Structures in Computer Science*, 25:1100–1115, 6 2015. (Cited on page 163.)
- [Voe06] Vladimir Voevodsky. A very short note on the homotopy λ -calculus. http://www.math.ias.edu/~vladimir/Site3/Univalent_Foundations_files/Hlambda_short_current.pdf, 2006. (Cited on page 3.)
- [Voe12] Vladimir Voevodsky. A universe polymorphic type system. <http://uf-ias-2012.wikispaces.com/file/view/Universe+polymorphic+type+system.pdf>, 2012. (Cited on pages 10 and 117.)
- [War08] Michael A. Warren. *Homotopy Theoretic Aspects of Constructive Type Theory*. PhD thesis, Carnegie Mellon University, 2008. (Cited on page 96.)
- [WR27] Alfred North Whitehead and Bertrand Russell. *Principia mathematica, 3 vols.* Cambridge University Press, Cambridge, 1910–1913; Second edition, 1925–1927. (Cited on pages 96 and 117.)

Index of symbols

$x : \equiv a$	definition, p. 19
$a \equiv b$	judgmental equality, p. 18
$a =_A b$	identity type, p. 45
$a = b$	identity type, p. 45
$x := b$	propositional equality by definition, p. 170
$\text{Id}_A(a, b)$	identity type, p. 45
$a =_p^p b$	dependent path type, p. 170
$a \neq b$	disequality, p. 51
refl_x	reflexivity path at x , p. 45
p^{-1}	path reversal, p. 58
$p \cdot q$	path concatenation, p. 59
$p \cdot_l r$	left whiskering, p. 64
$r \cdot_r q$	right whiskering, p. 64
$r \star s$	horizontal concatenation of 2-paths, p. 64
$g \circ f$	composite of functions, p. 53
$g \circ f$	composite of morphisms in a precategory, p. ??
f^{-1}	quasi-inverse of an equivalence, p. 73
f^{-1}	inverse of an isomorphism in a precategory, p. ??
0	empty type, p. 32
1	unit type, p. 26
\star	canonical inhabitant of 1 , p. 26
2	type of booleans, p. 33
$1_2, 0_2$	constructors of 2 , p. 33
$0_I, 1_I$	point constructors of the interval I , p. 173
AC	axiom of choice, p. 110
AC_∞	“type-theoretic axiom of choice”, p. 94
$\text{acc}(a)$	accessibility predicate, p. ??
$P \wedge Q$	logical conjunction (“and”), p. 109
$\text{ap}_f(p) \text{ or } f(p)$	application of $f : A \rightarrow B$ to $p : x =_A y$, p. 66
$\text{apd}_f(p)$	application of $f : \prod_{(a:A)} B(a)$ to $p : x =_A y$, p. 68
$\text{apd}_f^2(p)$	two-dimensional dependent ap, p. 176
$x \# y$	apartness of real numbers, p. ??
base	basepoint of S^1 , p. 167

base	basepoint of \mathbb{S}^2 , p. 168 and p. 175
$\text{biinv}(f)$	proposition that f is bi-invertible, p. 128
$x \sim y$	bisimulation, p. ??
$-$	blank used for implicit λ -abstractions, p. 21
\mathcal{C}	type of Cauchy approximations, p. ??
Card	type of cardinal numbers, p. ??
$\circ A$	reflector or modality applied to A , p. 230 and p. 232
$\text{cocone}_X(Y)$	type of cocones, p. 183
code	family of codes for paths, p. 86, p. ??, p. ??
$A \setminus B$	subset complement, p. 109
$\text{cons}(x, \ell)$	concatenation constructor for lists, p. 139 and p. 193
contr_x	path to the center of contraction, p. 114
$\mathcal{F} \triangleleft (J, \mathcal{G})$	inductive cover, p. ??
$\text{isCut}(L, U)$	the property of being a Dedekind cut, p. ??
$\{L \mid R\}$	cut defining a surreal number, p. ??
X^\dagger	morphism reversal in a \ddagger -category, p. ??
decode	decoding function for paths, p. 86, p. ??, p. ??
encode	encoding function for paths, p. 86, p. ??, p. ??
η_A^\circlearrowright or η_A	the function $A \rightarrow \circ A$, p. 230 and p. 232
$A \twoheadrightarrow B$	epimorphism or surjection
$\text{eq}_{\text{No}}(x, y)$	path constructor of the surreals, p. ??
$\text{eq}_{\mathbb{R}_c}(u, v)$	path constructor of the Cauchy reals, p. ??
$a \sim b$	an equivalence relation, p. 188
$X \simeq Y$	type of equivalences, p. 73
$\text{Equiv}(X, Y)$	type of equivalences (same as $X \simeq Y$)
$A \simeq B$	type of equivalences of categories, p. ??
$P \Leftrightarrow Q$	logical equivalence, p. 109
$\exists(x : A). B(x)$	logical notation for mere existential, p. 109
$\text{ext}(f)$	extension of $f : A \rightarrow B$ along η_A , p. 213
\perp	logical falsity, p. 109
$\text{fib}_f(b)$	fiber of $f : A \rightarrow B$ at $b : B$, p. 126
$\text{Fin}(n)$	standard finite type, p. 24
$\forall(x : A). B(x)$	logical notation for dependent function type, p. 109
funext	function extensionality, p. 80
$A \rightarrow B$	function type, p. 21
B^A	functor precategory, p. ??
glue	path constructor of $A \sqcup^C B$, p. 181
happly	function making a path of functions into a homotopy, p. 80
$\text{hom}_A(a, b)$	hom-set in a precategory, p. ??
$f \sim g$	homotopy between functions, p. 71
I	the interval type, p. 173
id_A	the identity function of A , p. 25

1_a	identity morphism in a precategory, p. ??
idtoeqv	function $(A = B) \rightarrow (A \simeq B)$ which univalence inverts, p. 83
idtoiso	function $(a = b) \rightarrow (a \cong b)$ in a precategory, p. ??
$\text{im}(f)$	image of map f , p. 226
$\text{im}_n(f)$	n -image of map f , p. 226
$P \Rightarrow Q$	logical implication (“implies”), p. 109
$a \in P$	membership in a subset or subtype, p. 106
$x \in v$	membership in the cumulative hierarchy, p. ??
$x \widetilde{\in} v$	resized membership, p. ??
ind_0	induction for $\mathbf{0}$, p. 33,
ind_1	induction for $\mathbf{1}$, p. 29,
ind_2	induction for $\mathbf{2}$, p. 34,
$\text{ind}_{\mathbb{N}}$	induction for \mathbb{N} , p. 37, and
$\text{ind}_{=A}$	path induction for $=_A$, p. 47,
$\text{ind}'_{=A}$	based path induction for $=_A$, p. 48,
$\text{ind}_{A \times B}$	induction for $A \times B$, p. 28,
$\text{ind}_{\sum_{(x:A)} B(x)}$	induction for $\sum_{(x:A)} B$, p. 30,
ind_{A+B}	induction for $A + B$, p. 33,
$\text{ind}_W_{(x:A)} B(x)$	induction for $W_{(x:A)} B$, p. 156
A/a	initial segment of an ordinal, p. ??
$\text{inj}(A, B)$	type of injections, p. ??
inl	first injection into a coproduct, p. 32
inr	second injection into a coproduct, p. 32
$A \cap B$	intersection of subsets, p. 109, classes, p. ??, or intervals, p. ??
$\text{isContr}(A)$	proposition that A is contractible, p. 114
$\text{isequiv}(f)$	proposition that f is an equivalence, p. 73, p. 121, and p. 129
$\text{ishae}(f)$	proposition that f is a half-adjoint equivalence, p. 124
$a \cong b$	type of isomorphisms in a (pre)category, p. ??
$A \cong B$	type of isomorphisms between precategories, p. ??
$A \cong B$	type of isomorphisms between sets, p. 73
$a \cong^\dagger b$	type of unitary isomorphisms, p. ??
isotoid	inverse of idtoiso in a category, p. ??
$\text{is-}n\text{-type}(X)$	proposition that X is an n -type, p. 205
$\text{isProp}(A)$	proposition that A is a mere proposition, p. 103
$\text{isSet}(A)$	proposition that A is a set, p. 99
$A * B$	join of A and B , p. 184
$\ker(f)$	kernel of a map of pointed sets, p. ??
$\lambda x. b(x)$	λ -abstraction, p. 25
$\text{lcoh}_f(g, \eta)$	type of left adjoint coherence data, p. 127
LEM	law of excluded middle, p. 104
LEM_∞	inconsistent propositions-as-types LEM, p. 102 and p. 105
$x < y$	strict inequality on natural numbers, p. 43, ordinals, p. ??, Cauchy reals, p. ??, surreals, p. ??, etc.

$x \leq y$	non-strict inequality on natural numbers, p. 43, Cauchy reals, p. ??, surreals, p. ??, etc.
\preceq, \prec	recursive versions of \leq and $<$ for surreals, p. ??
$\triangleleft, \triangleleft, \sqsubseteq, \sqsubset$	orderings on codomain of No-recursion, p. ??
$\lim(x)$	limit of a Cauchy approximation, p. ??
$\text{linv}(f)$	type of left inverses to f , p. 126
$\text{List}(X)$	type of lists of elements of X , p. 139 and p. 193
loop	path constructor of S^1 , p. 167
$\text{Map}_*(A, B)$	type of based maps, p. 178
$x \mapsto b$	alternative notation for λ -abstraction, p. 21
$\max(x, y)$	maximum in some ordering, e.g. p. ?? and p. ??
$\text{merid}(a)$	meridian of ΣA at $a : A$, p. 176
$\min(x, y)$	minimum in some ordering, e.g. p. ?? and p. ??
$A \rightarrowtail B$	monomorphism or embedding
\mathbb{N}	type of natural numbers, p. 35
\mathbb{N}	north pole of ΣA , p. 176
$\mathbb{N}^W, 0^W, \text{succ}^W$	natural numbers encoded as a W -type, p. 144
$\mathbb{N}\text{Alg}$	type of \mathbb{N} -algebras, p. 147
$\mathbb{N}\text{Hom}(C, D)$	type of \mathbb{N} -homomorphisms, p. 147
nil	empty list, p. 139 and p. 193
No	type of surreal numbers, p. ??
$\neg P$	logical negation (“not”), p. 109
$n\text{-Type}, n\text{-Type}_{\mathcal{U}}$	universe of n -types, p. 208
$\Omega(A, a), \Omega A$	loop space of a pointed type, p. 65
$\Omega^k(A, a), \Omega^k A$	iterated loop space, p. 65
A^{op}	opposite precategory, p. ??
$P \vee Q$	logical disjunction (“or”), p. 109
Ord	type of ordinal numbers, p. ??
(a, b)	(dependent) pair, p. 26 and p. 29
$\text{pair}^=$	constructor for $=_{A \times B}$, p. 76
$\pi_n(A)$	n^{th} homotopy group of A , p. 192 and p. ??
$\mathcal{P}(A)$	power set, p. 107
$\mathcal{P}_+(A)$	merely-inhabited power set, p. ??
pred	predecessor function $\mathbb{Z} \rightarrow \mathbb{Z}$, p. ??
$A \times B$	cartesian product type, p. 26
$\prod_{(x:A)} B(x)$	dependent function type, p. 24
$\text{pr}_1(t)$	the first projection from a pair, p. 27 and p. 30
$\text{pr}_2(t)$	the second projection from a pair, p. 27 and p. 30
$\text{Prop}, \text{Prop}_{\mathcal{U}}$	universe of mere propositions, p. 106
$A \times_C B$	pullback of A and B over C , p. 95
$A \sqcup^C B$	pushout of A and B under C , p. 181
\mathbb{Q}	type of rational numbers, p. ??
\mathbb{Q}_+	type of positive rational numbers, p. ??

$\text{qinv}(f)$	type of quasi-inverses to f , p. 72
A/R	quotient of a set by an equivalence relation, p. 187
$A \mathbin{\text{/\hspace{-0.1cm}/}} R$	alternative definition of quotient, p. 189
\mathbb{R}	type of real numbers (either), p. ??
\mathbb{R}_c	type of Cauchy real numbers, p. ??
\mathbb{R}_d	type of Dedekind real numbers, p. ??
$\text{rat}(q)$	rational number regarded as a Cauchy real, p. ??
$\text{rcoh}_f(g, \epsilon)$	type of right adjoint coherence data, p. 127
rec_0	recursor for 0 , p. 33
rec_1	recursor for 1 , p. 28
rec_2	recursor for 2 , p. 33
$\text{rec}_{\mathbb{N}}$	recursor for \mathbb{N} , p. 36
$\text{rec}_{A \times B}$	recursor for $A \times B$, p. 27
$\text{rec}_{\sum_{(x:A)} B(x)}$	recursor for $\sum_{(x:A)} B(x)$, p. 30
rec_{A+B}	recursor for $A + B$, p. 32
$\text{rec}_{W_{(x:A)} B(x)}$	recursor for $W_{(x:A)} B(x)$, p. 146
rinv	type of right inverses to f , p. 126
S	south pole of ΣA , p. 176
S^n	n -dimensional sphere, p. 174
seg	path constructor of the interval I , p. 173
$\text{Set}, \text{Set}_{\mathcal{U}}$	universe of sets, p. 106
Set	category of sets, p. ??
$\text{set}(A, f)$	constructor of the cumulative hierarchy, p. ??
$x \sim_{\epsilon} y$	relation of ϵ -closeness for \mathbb{R}_c , p. ??
$x \approx_{\epsilon} y$	recursive version of \sim_{ϵ} , p. ??
\sim_{ϵ} or $\curvearrowleft_{\epsilon}$	closeness relations on codomain of \mathbb{R}_c -recursion, p. ??
$A \wedge B$	smash product of A and B , p. 184
$\{ x : A \mid P(x) \}$	subset type, p. 106
$\{ f(x) \mid P(x) \}$	image of a subset, p. ??
$B \subseteq C$	containment of subset types, p. 106
$(q, r) \subseteq (s, t)$	inclusion of intervals, p. ??
succ	successor function $\mathbb{N} \rightarrow \mathbb{N}$, p. 35
succ	successor function $\mathbb{Z} \rightarrow \mathbb{Z}$, p. ??
$A + B$	coproduct type, p. 32
$\sum_{(x:A)} B(x)$	dependent pair type, p. 29
$\text{sup}(a, f)$	constructor for W -type, p. 144
surf	2-path constructor of \mathbb{S}^2 , p. 168 and p. 175
ΣA	suspension of A , p. 176
$\text{total}(f)$	induced map on total spaces, p. 132
$p_*(u)$	transport of $u : P(x)$ along $p : x = y$, p. 67
$\text{transport}^P(p, u)$	transport of $u : P(x)$ along $p : x = y$, p. 67
$\text{transport}^2(X, Y)$	two-dimensional transport, p. 175

$\text{transportconst}_Y^X(Z)$	transporting in a constant family, p. 69
$\ A\ _n$	n -truncation of A , p. 212
$ a _n^A, a _n$	image of $a : A$ in $\ A\ _n$, p. 212
$\ A\ $	propositional truncation of A , p. 108 and p. 185
$ a $	image of $a : A$ in $\ A\ $, p. 108 and p. 185
\top	logical truth, p. 109
$_$	an unnamed object or variable
$A \cup B$	union of subsets, p. 109
$\text{uniq}_{A \times B}$	uniqueness principle for the product $A \times B$, p. 28
uniq_1	uniqueness principle for 1 , p. 29
\mathcal{U}	universe type, p. 23
$\mathcal{U}_\circlearrowleft$	universe of modal types, p. 232
\mathcal{U}_\bullet	universe of pointed types, p. 65
ua	inverse to idtoeqv from univalence, p. 83
V	cumulative hierarchy, p. ??
$\text{WAlg}(A, B)$	type of w -algebras, p. 148
$\text{WHom}_{A,B}(C, D)$	type of W -homomorphisms, p. 148
$\text{W}_{(x:A)}B(x)$	W -type (inductive type), p. 144
$A \vee B$	wedge of A and B , p. 184
\mathbf{y}	Yoneda embedding, p. ??
\mathbb{Z}	type of integers, p. 189

Index

○-connected function, 233

○-truncated function, 233

∞-functor, 58

∞-group, 192

∞-groupoid, 3, 56, 58, 96, 147, 160, 168
fundamental, 56

structure of a type, 58–65

(∞, 1)-category, 97, 137, 147, 213

(∞, 1)-topos, 11, 134, 198, 233, 234

1-type, 100

2-dimensional path, *see* path, 2-

2-out-of-3 property, 131

2-out-of-6 property, 138

2-path, *see* path, 2-

3-dimensional path, *see* path, 3-

3-path, *see* path, 3-

abelian group, *see* group, abelian
abstraction

λ-, *see* λ-abstraction

abuse

of language, 113, 114

of notation, 4, 90

acceptance, 194

Ackermann function, 54

action

of a dependent function on a path, 68

of a function on a path, 66

addition

of natural numbers, 36

adjective, 113

adjoining a disjoint basepoint, 178

adjoint

equivalence, 72, 137

of types, half, 124–127

functor, 56, 94

functor theorem, 194

adjunction, *see* adjoint functor

adverb, 113, 233

AGDA, *see* proof assistant

algebra

2-cell, 149

colimits of, 195

for a polynomial functor, 148

for an endofunctor, 148

free, 193

initial, 155n, *see* homotopy-initial

ℕ-, 147

W-, 148

algorithm, 5, 7, 8, 19, 42, 57, 193

α-conversion, 22n

amalgamated free product, 195

analytic mathematics, 57

anger, 101–103, 210

application

of dependent function, 25

of dependent function to a path, 68

of function, 21

of function to a path, 66

of hypothesis or theorem, 41

arity, 144

associativity, 194

in a group, 192

in a monoid, 192

of addition

of natural numbers, 37

of function composition, 53

of function types, 23

of list concatenation, 193

of path concatenation, 61

coherence of, 63

of semigroup operation, 90

of Σ-types, 97

assumption, 19–20

attaching map, 179, 180

automorphism

fixed-point-free, 9, 102

of 2, nonidentity, 202, 204

of S^1 , 175

axiom

double negation, 104

excluded middle, *see* excluded middle

function extensionality, *see* function extensionality

of choice, 8, 9, 110–111, 118

$\text{AC}_{n,m}$, 235

- n*-connected, 235
- type-theoretic, 31, 94, 96, 101
 - unique, *see* unique choice
- of reducibility, 117
- propositional resizing, *see* propositional resizing
- Streicher’s Axiom K, 52, 209, 233
 - generalization to *n*-types, 210
- univalence, *see* univalence axiom
- unstable octahedral, 138
- versus rules, 20, 74, 170n
- axiomatic freedom, 42
- based map, 178
- basepoint, 65, 223
 - adjoining a disjoint, 178
- β -conversion, *see* β -reduction
- β -reduction, 22n, 26n
- bi-invertible function, 128
- bijection, 73, 130
- binding structure, 22
- bit, 103
- boolean
 - type of, *see* type of booleans
- bound variable, *see* variable, bound
- bounded
 - totally, *see* totally bounded
- Bourbaki, 96
- bracket type, *see* truncation, propositional
- canonicity, 10, 20
- capture, of a variable, 22
- carrier, 32, 90
- cartesian product, *see* type, product
- case analysis, 32, 140
- category
 - ($\infty, 1$)-, *see* ($\infty, 1$)-category
 - center of, 122
- caves, walls of, 142
- cell complex, 179–181
- center
 - of a category, 122
 - of contraction, 114
- choice operator, 102
- circle type, *see* type, circle
- classical
 - homotopy theory, 55–56
 - logic, *see* logic
 - mathematics, *see* mathematics, classical
- classifier
 - object, 134, 137
- closed
 - modality, 236
- cocone, 183, 217
- codes, *see* encode-decode method
- codomain, of a function, 21
- coequalizer, 197
- coercion, universe-raising, 52
- cofiber of a function, 184
- coherence, 61, 63, 124, 127, 151
- colimit
 - of sets, 186
 - of types, 95, 181, 234, 236
- commutative
 - group, *see* group, abelian
 - square, 97
- comonad, 234
- complement, of a subset, 109
- component, of a pair, *see* projection
- composition
 - of functions, 53
 - of paths, 59
 - horizontal, 64
- computation rule, 26
 - for coproduct type, 33
 - for dependent function types, 25
 - for dependent pair type, 30
 - for function types, 22
 - for higher inductive types, 169–170
 - for identity types, 47
 - for inductive types, 155
 - for natural numbers, 35, 37, 140
 - for product types, 27
 - for S^1 , 169, 171
 - for type of booleans, 140
 - for W-types, 146
 - propositional, 26, 150, 151, 169–170, 181
 - for identities between functions, 81
 - for identities between pairs, 76
 - for univalence, 83
- computational effect, 234
- computer proof assistant, *see* proof assistant
- concatenation of paths, 59
- cone
 - of a function, 184
 - of a sphere, 180
- conjunction, 39, 109
- connected
 - categorically, 237
 - function, *see* function, *n*-connected
 - type, 221
- consistency, 10, 44n, 155
- constant
 - function, 21
 - type family, 24
- constructive

- logic, *see* logic
- mathematics, *see* mathematics, constructive
- constructivity, 10
- constructor, 26, 155
 - path, 167
 - point, 167
- containment
 - of subsets, 106
- context, 19
- "continuity" of functions in type theory, 3, 66, 68, 71, 80, 99, 102, 115, 210
- continuous map, *see* function, continuous
- contractible
 - function, 128–129
 - type, 114–116
- contradiction, 40
- contravariant functor, 154
- conversion
 - α -, *see* α -conversion
 - β -, *see* β -reduction
 - η -, *see* η -expansion
- coproduct, *see* type, coproduct
- CoQ, *see* proof assistant
- corollary, 17n
- countable axiom of choice, *see* axiom of choice, countable
- covariant functor, 153
- cumulative
 - universes, 23
- currying, 23
- CW complex, 5, 179–181
- de Morgan's laws, 40–42
- decidable
 - definitional equality, 19
 - equality, 43, 106, 195, 210
 - subset, 35
 - type, 105
 - type family, 106
- deductive system, 17
- definition
 - by pattern matching, 38–39, 156
 - inductive, 139, *see* type, inductive
 - of adverbs, 113
 - of function, direct, 21, 24
- definitional equality, *see* equality, definitional
- denial, 104–106, 110–111
- dependent
 - function, *see* function, dependent
 - path, *see* path, dependent
 - type, *see* type, family of
- dependent eliminator, *see* induction principle
- diagram, 71, 97, 234, 236
- dimension
 - of path constructors, 168
 - of paths, 56
- disc, 179, 180
- discrete
 - space, 6, 8, 11, 99
- disequality, 51
- disjoint
 - basepoint, 178
 - sum, *see* type, coproduct
 - union, *see* type, coproduct
- disjunction, 39, 109
- domain
 - of a constructor, 153
 - of a function, 21
- double negation, law of, 101, 104
- dummy variable, *see* variable, bound
- Eckmann–Hilton argument, 64, 192
- effective
 - procedure, 42
- element, 18
- Elementary Theory of the Category of Sets, 6
- elimination rule, *see* eliminator
- eliminator, 26
 - of inductive type
 - dependent, *see* induction principle
 - non-dependent, *see* recursion principle
 - embedding, *see* function, embedding
 - empty type, *see* type, empty
 - encode-decode method, 88, 86–90
 - end point of a path, 56
- endofunctor
 - algebra for, 163
 - polynomial, 148, 155n
- equality
 - decidable, *see* decidable equality
 - definitional, 18, 169
 - heterogeneous, 170
 - judgmental, 18, 169
 - merely decidable, 210
 - propositional, 18, 45
 - reflexivity of, 60
 - symmetry of, 58, 60
 - transitivity of, 59, 60
 - type, *see* type, identity
- equals may be substituted for equals, 45
- equipped with, 31
- equivalence, 72–74, 83–84, 129
 - as bi-invertible function, 128
 - as contractible function, 128–129
 - class, 188
 - fiberwise, 133

-
- half adjoint, 124–127
 - induction, 162
 - logical, 44
 - properties of, 129, 131
 - relation, *see* relation, equivalence
 - η -conversion, *see* η -expansion
 - η -expansion, 22n, 26n
 - Euclid of Alexandria, 56
 - evaluation, *see* application, of a function
 - evidence, of the truth of a proposition, 18, 39
 - ex falso quodlibet*, 33
 - excluded middle, 8, 9, 42, 102, 104, 210
 - LEM_{n,m}, 235
 - existential quantifier, *see* quantifier, existential
 - expansion, η -, *see* η -expansion
 - exponential ideal, 230
 - extensional
 - type theory, 51, 95
 - extensionality, of functions, *see* function extensionality
 - extraction of algorithms, 7, 8
 - f*-local type, 236
 - factorization
 - stability under pullback, 229
 - system, orthogonal, *see* orthogonal factorization system
 - false, 39, 40, 109
 - family
 - of types, *see* type, family of
 - Feit–Thompson theorem, 6
 - fiber, 126
 - fiberwise, 68
 - equivalence, 133
 - map, *see* fiberwise transformation
 - n -connected family of functions, 223
 - transformation, 132–133, 224
 - fibration, 67, 78, 132
 - finite
 - lists, type of, 193
 - sets, family of, 24, 25, 54
 - first-order
 - logic, 17, 112
 - fixed-point property, 164
 - flattening lemma, 197
 - formation rule, 26
 - foundations, 1
 - foundations, univalent, 1
 - four-color theorem, 7
 - free
 - algebraic structure, 193
 - generation of an inductive type, 140, 160, 167
 - group, *see* group, free
 - monoid, *see* monoid, free
 - product, 196
 - amalgamated, 195
 - Frege, 162
 - function, 21–23, 66
 - connected, 233
 - truncated, 233
 - Ackermann, 54
 - application, 21
 - application to a path of, 66
 - bi-invertible, 121, 128
 - bijective, *see* bijection
 - codomain of, 21
 - composition, 53
 - constant, 21
 - “continuity” of, *see* “continuity”
 - continuous
 - in classical homotopy theory, 2
 - contractible, 128–129
 - currying of, 23
 - dependent, 24–25, 68–70
 - application, 25
 - application to a path of, 68
 - domain of, 21
 - embedding, 130, 206, 225
 - fiber of, *see* fiber
 - fiberwise, *see* fiberwise transformation
 - “functoriality” of, *see* “functoriality”
 - idempotent, 189
 - identity, 25, 72, 83
 - injective, 130, 225
 - λ -abstraction, *see* λ -abstraction
 - left invertible, 126
 - n -connected, 221, 221
 - n -image of, 226
 - n -truncated, 225
 - polymorphic, 25
 - projection, *see* projection
 - quasi-inverse of, *see* quasi-inverse
 - retraction, 115, 130
 - right invertible, 126
 - section, 115
 - split surjective, 130
 - surjective, 130, 221
 - function extensionality, 52, 74, 80, 96, 115, 136
 - non-dependent, 138
 - proof from interval type, 174
 - proof from univalence, 135
 - weak, 135
 - function type, *see* type, function
 - functional relation, 21
 - functor
 - contravariant, 154

- covariant, 153
 polynomial, *see* endofunctor, polynomial
 “functoriality” of functions in type theory, 66, 71,
 80, 99, 102, 115, 210
 fundamental
 ∞ -groupoid, 56
 group, 55, **192**
- game
 deductive system as, 17
 generation
 of a type, inductive, 46–49, 139–141, 160, 167–
 169
 of an ∞ -groupoid, 168
 generator
 of a group, 193
 of an inductive type, *see* constructor
 geometric realization, 56
 geometry, synthetic, 56
 globular operad, 63
 graph, 97, **234**
 with composition, **236**
 Grothendieck construction, 197
 group, **192**
 abelian, 65, 123, 192
 free, 193–195
 fundamental, *see* fundamental group
 homomorphism, **194**
 homotopy, *see* homotopy group
 groupoid
 ∞ -, *see* ∞ -groupoid
 higher, 63
- h-initial, *see* homotopy-initial
 h-level, *see* n -type
 h-proposition, *see* mere proposition
 half adjoint equivalence, **124–127**
 Haskell, 234
 Hedberg’s theorem, 117, 210
 heterogeneous equality, 170
 hierarchy
 of n -types, *see* n -type
 of universes, *see* type, universe
 higher category theory, 55–56
 higher groupoid, *see* ∞ -groupoid
 higher inductive type, *see* type, higher inductive
 higher topos, *see* $(\infty, 1)$ -topos
 homomorphism
 group, **194**
 monoid, **193**
 \mathbb{N} -, **147**
 of algebras for a functor, **148**
 semigroup, 92
- W-, **148**
 homotopy, **70–72**, 80–82
 colimit, *see* colimit of types
 equivalence, *see* equivalence
 topological, 2
 fiber, *see* fiber
 group, 192
 hypothesis, 56
 induction, 162
 limit, *see* limit of types
 n -type, *see* n -type
 theory, classical, *see* classical homotopy theory
 topological, 2, 55
 type, 3
 homotopy-inductive type, 150
 homotopy-initial
 algebra for a functor, **149**
 \mathbb{N} -algebra, **147**
 W-algebra, **149**
 horizontal composition
 of paths, 64
 hub and spoke, **180–181**, 212
 hypothesis, **20**, 41, 42
 homotopy, 56
 inductive, 37
- idempotent
 function, **189**
 modality, 233
 identification, **45**
 identity, 4
 function, 25, 72, 83
 modality, 233
 system, **161–162**
 at a point, **160–161**
 type, *see* type, identity
 image, **226**
 n -image, **226**
 stability under pullback, 229
 implementation, *see* proof assistant
 implication, 39, 40, **109**
 impredicative
 encoding of a W-type, 165
 quotient, 188
 truncation, 117
 impredicativity, *see* mathematics, predicative
 for mere propositions, *see* propositional resizing
 inaccessible cardinal, 10
 inclusion
 of subsets, **106**
 index of an inductive definition, **157**
 indiscernability of identicals, 45, 67

- induction principle, 29, 140
 - for a modality, 232
 - for an inductive type, 156
 - for connected maps, 222
 - for coproduct, 33
 - for dependent pair type, 30
 - for empty type, 33
 - for equivalences, 162
 - for homotopies, 162
 - for identity type, 46–49, 57
 - based, 47
 - for integers, 190
 - for interval type, 173
 - for natural numbers, 36
 - for product, 29
 - for S^1 , 170, 174
 - for S^2 , 176
 - for suspension, 176
 - for torus, 180
 - for truncation, 118, 185, 212
 - for type of booleans, 34
 - for type of vectors, 157
 - for W-types, 145
- inductive
 - definition, 139, *see* type, inductive hypothesis, 37
 - predicate, 157
 - type, *see* type, inductive
 - higher, *see* type, higher inductive type family, 157
 - inductive-inductive type, 158
 - inductive-recursive type, 159
 - inequality, *see* order
 - infinitary
 - algebraic theory, 196
 - informal type theory, 6–7
 - inhabited type, 44, 103, 105
 - merely, 113
 - initial
 - algebra characterization of inductive types, *see* homotopy-initial field, 196
 - type, *see* type, empty
 - injection, *see* function, injective
 - injective function, *see* function, injective
 - integers, 189
 - induction principle for, 190
 - intensional type theory, 51, 95
 - interchange law, 65
 - intersection
 - of subsets, 109
 - interval
 - topological unit, 3
 - type, *see* type, interval
- introduction rule, 26
- intuitionistic logic, *see* logic
- inverse
 - in a group, 192
 - left, 126
 - of path, 58
 - right, 126
- isomorphism
 - natural, 70
 - of sets, 73, 130
 - semigroup, 92
 - transfer across, 143
- iterated loop space, 64
- iterator
 - for natural numbers, 53
- J*, *see* induction principle for identity type
- join
 - of types, 184
- judgment, 17, 18
- judgmental equality, 18, 74, 141, 169
- k*-morphism, 56
- Kan complex, 4, 10
- Klein bottle, 180
- λ -abstraction, 21, 23, 25, 38, 43
- λ -calculus, 2
- language, abuse of, *see* abuse of language
- law
 - de Morgan's, 40–42
 - of double negation, 104
 - of excluded middle, *see* excluded middle
- Lawvere, 6, 8, 51, 163, 164, 233
- lax colimit, 197, 198
- left
 - inverse, 126
 - invertible function, 126
- lemma, 17n
 - flattening, 197
- level, *see* universe level or *n*-type
- lifting
 - equivalences, 91
 - path, 67
- limit
 - of sets, 186
 - of types, 95, 97, 181
- list, *see* type of lists
- list type, *see* type, of lists
- logic
 - classical vs constructive, 41–42
 - constructive, 9
 - constructive vs classical, 9, 40, 101–106

- intuitionistic, 9
 of mere propositions, 103–104, 107–109, 112–114
 predicate, 42
 propositional, 39
 propositions as types, 39–44
 truncated, 112
 logical equivalence, 44
 logical notation, traditional, 109
 loop, 55, 63, 65
 constant, *see* path, constant
 dependent n --, 176, 203, 204
 n -, 65, 176, 205
 n -dimensional, *see* loop, n -loop space, 63, 65, 178, 192, 210
 iterated, 64, 65, 176, 179, 192, 211, 212
 n -fold, *see* loop space, iterated
- magma, 31, 43
 map, *see* function
 fiberwise, *see* fiberwise transformation
 of spans, 218
 mapping, *see* function
 mapping cone, *see* cone of a function
 Martin-Löf, 163
 matching, *see* pattern matching
 mathematics
 classical, 7, 8, 35, 41, 51, 94, 101, 103–105, 107, 110
 constructive, 7–10, 36, 192
 formalized, 2, 6, 6–7, 11, 19n, 60, 129, 158
 predicative, 107, 155
 proof-relevant, 20, 31, 44, 60, 72, 73, 198
 membership, 18
 mere proposition, 103–104, 106–109, 112–114
 mere relation, 187
 merely, 113, 233
 decidable equality, 210
 inhabited, 113
 meridian, 176, 180
 mistaken identity type, 239
 modal
 logic, 233
 operator, 233, 234, 236
 type, 232
 modality, 232, 230–234
 closed, 236
 identity, 233
 open, 236
 model category, 4
 monad, 203, 234
 monoid, 192, 192–195
 free, 193, 204
- homomorphism, 193
 morphism
 in an ∞ -groupoid, 56
 multiplication
 in a group, 192
 in a monoid, 192
 of natural numbers, 54
 mutual inductive type, 158
- \mathbb{N} -algebra, 147
 homotopy-initial (h-initial), 147
- n -connected
 axiom of choice, 235
 function, *see* function, n -connected
 type, *see* type, n -connected
- n -dimensional loop, *see* loop, n -
 n -dimensional path, *see* path, n -
 \mathbb{N} -homomorphism, 147
- n -image, 226
 n -loop, *see* loop, n -
 n -path, *see* path, n -
 n -sphere, *see* type, n -sphere
 n -truncated
 function, 225
 type, *see* n -type
 n -truncation, *see* truncation
 n -type, 8, 100–101, 205, 205–229
 definable in type theory, 117
- natural
 transformation, 122
- natural numbers, 35–37, 88–90, 140, 190
 as homotopy-initial algebra, 147
 encoded as a W-type, 144, 152
 encoded as List(1), 143
 isomorphic definition of, 141
 “naturality” of homotopies, 71
- negation, 40, 105
 negative
 type, 86
- non-dependent eliminator, *see* recursion principle
 notation, abuse of, *see* abuse of notation
 noun, 113
 nullary
 coproduct, *see* type, empty
 product, *see* type, unit
- number
 integers, 189
 natural, *see* natural numbers
 real, *see* real numbers
- NUPRL, *see* proof assistant
- object
 classifier, 134, 137

-
- subterminal, *see* mere proposition
 octahedral axiom, unstable, 138
 odd-order theorem, 6
 open
 modality, 236
 problem, 10–12, 97, 235, 236
 operad, 63
 operator
 choice, *see* choice operator
 induction, *see* induction principle
 modal, *see* modality
 orthogonal factorization system, 205, 225–229, 233
 pair
 dependent, 29
 ordered, 26
 paradox, 23, 107, 154
 parallel paths, 56
 parameter
 of an inductive definition, 157
 space, 19
 parentheses, 21, 23
 path, 45, 56, 58–65
 2-, 55, 56, 63
 2-dimensional, *see* path, 2-
 3-, 56, 56, 63
 3-dimensional, *see* path, 3-
 application of a dependent function to, 68
 application of a function to, 66
 composite, 59
 concatenation, 59
 n-fold, 191
 constant, 45, 48
 constructor, 167
 dependent, 69, 170
 in dependent function types, 82
 in function types, 82
 in identity types, 86
 end point of, 56
 induction, 46–49
 induction based, 47
 inverse, 58
 lifting, 67
 n-, 63, 97, 176, 203
 n-dimensional, *see* path, n-
 parallel, 56
 start point of, 56
 topological, 3, 55
 pattern matching, 38–39, 52, 156
 Peano, 162
 pentagon, Mac Lane, 56
 Π -type, *see* type, dependent function
 point
 constructor, 167
 of a type, 18
 pointed
 predicate, 160
 type, *see* type, pointed
 pointwise
 equality of functions, 80
 functionality, 53
 operations on functions, 81
 polarity, 86
 pole, 176
 polymorphic function, 25
 polynomial functor, *see* endofunctor, polynomial
 positive
 type, 86
 positivity, strict, *see* strict positivity
 Postnikov tower, 8, 212
 power set, 35, 107, 154, 188
 predecessor, 140, 145
 function, truncated, 170
 predicate
 inductive, 157
 logic, 42
 pointed, 160
 predicative mathematics, *see* mathematics, predicative
 presentation
 of a group, 196
 of a positive type by its constructors, 86
 of a space as a CW complex, 5
 of an ∞ -groupoid, 168
 prime number, 112
 primitive
 recursion, 35
 principle, *see* axiom
 product
 of types, *see* type, product
 programming, 2, 9, 23, 142, 234
 projection
 from cartesian product type, 27
 from dependent pair type, 30
 projective plane, 180
 proof, 18, 39–44
 assistant, 2, 7, 51, 203
 NUPRL, 117
 AGDA, 52
 COQ, 52, 96, 117
 by contradiction, 40, 42, 105
 proof-relevant mathematics, *see* mathematics, proof-relevant
 proposition
 as types, 7, 39–44, 101–103
 mere, *see* mere proposition

- propositional
 - equality, 18, 45
 - logic, 39
 - resizing, 107, 117, 118, 155
 - truncation, *see* truncation
 - uniqueness principle, *see* uniqueness principle, propositional
- pullback, 95, 97, 134, 182, 207
- purely, 113, 233, 234
- pushout, 182–185
 - in n -types, 218
 - of sets, 187
- quantifier, 42, 109
 - existential, 42, 108, 109
 - universal, 42, 108, 109
- quasi-inverse, 72, 122–124
- Quillen model category, 4
- quotient of sets, *see* set-quotient
- recurrence, 35, 140, 163
- recursion
 - primitive, 35
- recursion principle, 140
 - for a modality, 232
 - for an inductive type, 155
 - for cartesian product, 28
 - for coproduct, 32
 - for dependent pair type, 30
 - for empty type, 33
 - for interval type, 173
 - for natural numbers, 35
 - for S^1 , 169, 171
 - for S^2 , 175
 - for suspension, 176
 - for truncation, 108, 185, 213
 - for type of booleans, 33
- recursive call, 35
- recursor, *see* recursion principle
- reduced word in a free group, 195
- reduction
 - β -, *see* β -reduction
 - of a word in a free group, 195
- reflection rule, 95
- reflective
 - subcategory, 213
 - subuniverse, 230
- reflexivity
 - of a relation, 188, 209
 - of equality, 45
- relation
 - equivalence, 188
 - mere, 187
- reflexive, 188, 209
- symmetric, 188
- transitive, 188
- resizing, 117
 - propositional, *see* propositional resizing
- retract
 - of a function, 131–132, 221
 - of a type, 115, 137, 206
- retraction, 115, 130, 206
- right
 - inverse, 126
 - invertible function, 126
- rule, 17
 - elimination, *see* eliminator
 - formation, 26
 - introduction, 26
 - versus axioms, 20, 170n
- Russell, Bertrand, 2, 117
- scope, 21, 25, 29
- Scott, 164
- section, 115
 - of a type family, 68
- semigroup, 43, 90
 - structure, 90
- semiring, 54
- sequence, 57, 155
 - Cauchy, *see* Cauchy sequence
- set, 6, 99–101, 130, 186–187, 206, 208–210
- set theory
 - Zermelo–Fraenkel, 6, 17, 196
- set-pushout, 187
- set-quotient, 187–191, 203
- setoid, 117, 203
- Σ -type, *see* type, dependent pair
- signature
 - of an algebraic theory, 196
- simplicial
 - sets, 4
- simply connected type, 221
- singleton type, *see* type, singleton
- skeleton
 - of a CW-complex, 179
- small
 - type, 4, 24
- smash product, 184
- source
 - of a function, *see* domain
 - of a path constructor, 167, 179, 202
- space
 - topological, *see* topological space
- span, 182, 187, 217
- sphere type, *see* type, sphere

- split
 - surjection, *see* function, split surjective
 - spoke, *see* hub and spoke
 - squash type, *see* truncation, propositional
 - stability
 - and descent, 198
 - of images under pullback, 229
 - stages, five, of accepting constructive mathematics, 239
 - start point of a path, 56
 - strict
 - positivity, 155, 201
 - structure
 - semigroup, 90
 - subset, 106
 - subsingleton, *see* mere proposition
 - substitution, 20, 22
 - subterminal object, *see* mere proposition
 - subtype, 43, 106
 - subuniverse, reflective, 230
 - successor, 35, 90, 140, 144, 145
 - sum
 - dependent, *see* type, dependent pair
 - disjoint, *see* type, coproduct
 - of numbers, *see* addition
 - supremum
 - constructor of a W-type, 144
 - surjection, *see* function, surjective
 - split, *see* function, split surjective
 - surjective
 - function, *see* function, surjective
 - split, *see* function, split surjective
 - suspension, 176–179, 184
 - symmetry
 - of a relation, 188
 - of equality, 58
 - synthetic mathematics, 56
 - system, identity, *see* identity system
 - target
 - of a function, *see* codomain
 - of a path constructor, 167, 179, 202
 - term, 18
 - terminal
 - type, *see* type, unit
 - theorem, 17n
 - Feit–Thompson, 6
 - four-color, 7
 - Hedberg’s, 117, 210
 - odd-order, 6
 - theory
 - algebraic, 196
 - essentially algebraic, 196
 - Tierney, 233
 - together with, 31
 - topological
 - path, 3, 55
 - space, 2, 3, 55
 - topology
 - Lawvere-Tierney, 233
 - topos, 9, 117, 233, 234
 - higher, *see* $(\infty, 1)$ -topos
 - torus, 179, 181, 204
 - induction principle for, 180
 - total
 - space, 67, 78, 132
 - traditional logical notation, 109
 - transformation
 - fiberwise, *see* fiberwise transformation
 - transitivity
 - of a relation, 188
 - of equality, 59
 - transport, 67–70, 78, 83
 - in coproduct types, 88
 - in dependent function types, 81
 - in dependent pair types, 79
 - in function types, 81
 - in identity types, 85
 - in product types, 77
 - in unit type, 80
 - tree, well-founded, 144
 - true, 39, 109
 - truncation
 - n -truncation, 187, 212–217
 - propositional, 108–109, 185–186
 - set, 186–187
 - type
 - 2-sphere, 168, 175–176
 - bracket, *see* truncation, propositional
 - cartesian product, *see* type, product
 - circle, 167, 169–175, 177
 - coequalizer, 197
 - colimit, 95, 181
 - connected, 221
 - contractible, 114–116
 - coproduct, 32–33, 33, 34, 53, 86–88
 - decidable, 105
 - dependent, *see* type, family of
 - dependent function, 24–25, 80–82
 - dependent pair, 29–32, 77–79, 93
 - dependent sum, *see* type, dependent pair
 - empty, 32–33, 95, 144
 - equality, *see* type, identity
 - f -local, 236
 - family of, 24, 34, 67–70
 - constant, 24

- decidable, 106
- inductive, 157
- function, 21–23, 80–82
- higher inductive, 5, 167–204
- homotopy-inductive, 150
- identity, 45–51, 58–65, 84–86, 94
 - as inductive, 159
- inductive, 139–141, 153–159
 - generalizations, 156
- inductive-inductive, 158
- inductive-recursive, 159
- inhabited, *see* inhabited type
- interval, 173–174
- limit, 95, 181
- mistaken identity, 239
- modal, 232
- mutual inductive, 158
- n -connected, 221
- n -sphere, 176, 178, 179
- n -truncated, *see* n -type
- n -type, *see* n -type
- negative, 86
- of booleans, 33–35
- of lists, 139, 142, 144, 157, 163, 193
- of natural numbers, *see* natural numbers
- of vectors, 157
- Π -, *see* type, dependent function
- pointed, 65, 176
- positive, 86
- product, 26–29, 53, 75–77, 93
- pushout of, *see* pushout
- quotient, *see* set-quotient
- Σ -, *see* type, dependent pair
- simply connected, 221
- singleton, 48, 114
- small, 4, 24
- squash, *see* truncation, propositional
- subset, 106
- suspension of, *see* suspension
- truncation of, *see* truncation
- unit, 26–29, 33, 80, 95, 114, 143, 144
- universe, 23–24, 83–84, 101
 - cumulative, 23
 - level, *see* universe level
 - Russell-style, 52
 - Tarski-style, 52
 - univalent, 83
- W -, *see* W -type
- type theory, 2, 17
 - extensional, 51, 95
 - formal, 6–7
 - informal, 6–7, 20
 - intensional, 51, 95
- unintentional, 239
- typical ambiguity, 24
- UIP, *see* uniqueness of identity proofs
- unequal, 51
- unintentional type theory, 239
- union
 - disjoint, *see* type, coproduct of subsets, 109
- unique
 - choice, 111–112
 - factorization system, 225–229, 233
- uniqueness
 - of identity proofs, 52, 208
 - of identity types, 141
 - principle, 26, 52
 - for dependent function types, 25
 - for function types, 22, 52
 - for identities between functions, 81
 - for product types, 52
 - principle, propositional, 26
 - for a modality, 232
 - for dependent pair types, 79
 - for functions on a pushout, 183
 - for functions on a truncation, 213
 - for functions on \mathbb{N} , 141
 - for functions on the circle, 172
 - for functions on W -types, 146
 - for homotopy W -types, 151
 - for identities between pairs, 76
 - for product types, 28, 76
 - for univalence, 84
- unit
 - interval, 3
 - law for path concatenation, 61
 - of a group, 192
 - of a monoid, 192
 - type, *see* type, unit
- univalence axiom, 1, 4, 8, 74, 83, 91, 96, 101, 135, 143, 162
 - constructivity of, 10
- univalent universe, 83
- universal
 - property, 93–95
 - of W -type, 148
 - of a modality, 232
 - of cartesian product, 93
 - of coproduct, 97
 - of dependent pair type, 94, 198
 - of free group, 194
 - of identity type, 94
 - of natural numbers, 147
 - of pushout, 183, 187, 219

of S^1 , 172
of S^n , 179
of suspension, 178
of truncation, 186, 213
quantifier, *see* quantifier, universal
universal property, 144
universe, *see* type, universe
universe level, 24

value
of a function, 21
truth, 35

variable, 19, 21, 24, 38, 41, 43, 145, 156
bound, 22
captured, 22
dummy, 22
scope of, 21
type, 153, 157

vary along a path constructor, 171

vector, 157
induction principle for, 157

vertex of a cocone, 183

W-algebra, 148

W-homomorphism, 148

W-type, 144
as homotopy-initial algebra, 148
impredicative encoding of, 165

wedge, 184

whiskering, 64, 211

witness
to the truth of a proposition, 18, 39

Yoneda
lemma, 160

Zermelo-Fraenkel set theory, *see* set theory

zero, 35, 140, 144, 145

ZF, *see* set theory

ZFC, *see* set theory