

# THE RELATIVE CONSISTENCY OF THE AXIOM OF CHOICE MECHANIZED USING ISABELLE/ZF

LAWRENCE C. PAULSON

## Abstract

The proof of the relative consistency of the axiom of choice has been mechanized using Isabelle/ZF, building on a previous mechanization of the reflection theorem. The heavy reliance on metatheory in the original proof makes the formalization unusually long, and not entirely satisfactory: two parts of the proof do not fit together. It seems impossible to solve these problems without formalizing the metatheory. However, the present development follows a standard textbook, Kenneth Kunen's *Set theory: an introduction to independence proofs*, and could support the formalization of further material from that book. It also serves as an example of what to expect when deep mathematics is formalized.

## 1. Introduction

In 1940, Gödel [5] published his famous monograph proving that the axiom of choice (AC) and the generalized continuum hypothesis (GCH) are consistent with respect to the other axioms of set theory. This theorem addresses the first of Hilbert's celebrated list of mathematical problems. I have attempted to reproduce this work in Isabelle/ZF.

When so much mathematics has already been checked mechanically, what is the point of checking any more? Obviously, the theorem's significance makes it a challenge, as do its size and complexity, but the real challenge comes from its reliance on metamathematics. As I have previously noted [18], some theorems seem difficult to formalize even in their statements, let alone in their proofs. Gödel's work is not a single formal theorem. It consists of several different theorems which, taken collectively, can be seen as expressing the relative consistency of the axiom of choice. At the end of [5, Chapter VII], Gödel remarks that given a contradiction from the axioms of set theory augmented with the axiom of choice, a contradiction in basic set theory 'could actually be constructed' [5, p. 87]. This claim is crucial: logicians prefer consistency proofs to be constructive.

Gödel's idea [4, 6] is to define a very lean model, called **L**, of set theory. **L** contains just the sets that must exist because they can be defined by formulae. Then, prove that **L** satisfies the ZF axioms and the additional axiom that 'every set belongs to **L**', which is abbreviated to  $\mathbf{V} = \mathbf{L}$ . We now know that  $\mathbf{V} = \mathbf{L}$  is consistent with Zermelo–Fraenkel (ZF) set theory, and can assume that this axiom holds. (The conjunction of ZF and  $\mathbf{V} = \mathbf{L}$  is abbreviated as ZFL.) We conclude by proving that AC and GCH are theorems of ZFL, and therefore are also consistent with ZF.

A complication in Gödel's proof is its use of classes. Intuitively speaking, a *class* is a collection of sets that is defined by comprehension,  $\{x \mid \phi(x)\}$ . Every set  $A$  is trivially

---

Received 7 January 2003, revised 4 September 2003; published 13th October 2003.

2000 Mathematics Subject Classification 03E35, 03B35, 68T15

© 2003, Lawrence C. Paulson

a class, namely  $\{x \mid x \in A\}$ , but a *proper class* is too big to be a set. Formal set theories restrict the use of classes in order to eliminate the danger of paradoxes. Modern set theorists use Zermelo–Fraenkel (ZF) set theory, where classes exist only in the metalanguage. That is, the class  $\{x \mid \phi(x)\}$  is just an alternative notation for the formula  $\phi(x)$ , and  $a \in \{x \mid \phi(x)\}$  is just an alternative notation for  $\phi(a)$ . The universal class,  $\mathbf{V}$ , corresponds to the formula *True*. An ‘equation’ like  $\mathbf{V} = \bigcup_{\alpha \in \mathbf{ON}} V_\alpha$  stands for  $\forall x. \exists \alpha. \mathbf{ON}(\alpha) \wedge x \in V_\alpha$ . (Here,  $\mathbf{ON}$  denotes the class of ordinal numbers.) Gödel’s monograph [5] uses von Neumann–Bernays–Gödel (NBG) set theory, which allows quantification over classes but restricts their use in other ways. With either axiom system, classes immensely complicate the reasoning, making it highly syntactic.

Why did Gödel use classes? Working entirely with sets, he could have used essentially the same techniques to prove that if  $M$  is a model of ZF, then there exists a model  $L(M)$  of ZFC. (ZFC refers to the ZF axioms plus AC.) Therefore, if ZFC has no models, then neither does ZF. But with this approach, he can no longer claim that if he had a contradiction in ZFC, then a contradiction in ZF ‘could actually be constructed’.<sup>1</sup> For the sake of this remark, which is not part of any theorem statement, Gödel chose a more difficult route. Classes create more difficulties for formal proof checkers (which have to face foundational issues squarely) than they do for mathematicians writing in English.

The proof uses metatheoretic reasoning extensively. Gödel writes [5, p. 34]:

*However, the only purpose of these general metamathematical considerations is to show how the proofs for theorems of a certain kind can be accomplished by a general method. And, since applications to only a finite number of instances are necessary . . . , the general metamathematical considerations could be left out entirely, if one took the trouble to carry out the proofs separately for any instance.*

I decided to take the trouble, not using metatheory, but relying instead on a mechanical theorem prover.

This paper describes the Isabelle/ZF proofs. It indicates the underlying mathematical ideas, and sometimes discusses practical issues such as proof length or the machine resources used. It necessarily omits much material that would be too long or too repetitious. The paper concerns how existing mathematics is formalized; it contains no original mathematics.

**Overview.** The paper begins by outlining Gödel’s relative consistency proof (Section 2). After a brief overview of Isabelle/ZF, the paper describes the strategy guiding the formalization (Section 3) and presents some elementary absoluteness proofs (Section 4). It then discusses relativization issues involving well-founded recursion (Section 5). Turning away from absoluteness, the paper proceeds to describe the formalization of the constructible universe and the proof that  $\mathbf{L}$  satisfies the ZF axioms (Section 6); then it describes how the reflection theorem is used to prove that  $\mathbf{L}$  satisfies the separation axiom (Section 7). Absoluteness again takes centre stage as the paper presents the relativization of two essential datatypes (Section 8) and finally presents the absoluteness of  $\mathbf{L}$  itself (Section 9). Finally, the paper presents the Isabelle proof that AC holds in  $\mathbf{L}$  (Section 10), and offers some conclusions (Section 11). Appendix A, which was generated by Isabelle/ZF, presents the full mechanical development.

<sup>1</sup>Consistency proofs using sets in ZF can be constructive, provided that we take care to keep track of the number of instances of the separation axiom that are required. We can prove in ZF – using the reflection theorem – that every finite conjunction of the ZF axioms has a model. Therefore, the absence of a model yields a contradiction in ZF. Proofs using *forcing* are often formalized using this technique. I am indebted to Kunen for this observation.

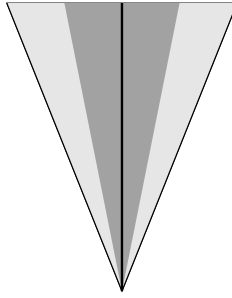


Figure 1: The constructible universe,  $\mathbf{L}$ .

## 2. Proof outline

Recall that Gödel's idea is to define a lean model of set theory, the class  $\mathbf{L}$  of the *constructible sets*. Figure 1 shows  $\mathbf{L}$  (shaded) as a subclass of the universe,  $\mathbf{V}$ . The vertical line represents the class  $\mathbf{ON}$  of the ordinals.

Gödel's proof involves four main tasks:

1. defining the class  $\mathbf{L}$  within ZF;
2. proving that  $\mathbf{L}$  satisfies the ZF axioms;
3. proving that  $\mathbf{L}$  satisfies  $\mathbf{V} = \mathbf{L}$ ;
4. proving that  $\mathbf{V} = \mathbf{L}$  implies the axiom of choice.

As we shall see, Isabelle is well suited to completing the first and last parts. Both the definition of  $\mathbf{L}$  and the proof of the axiom of choice are straightforward exercises in mechanized set theory. The second task cannot quite be completed: all of the ZF axioms can be verified, apart from *separation* – it is an axiom scheme, and each instance requires its own proof. As for the third task, the Isabelle proof that  $\mathbf{L}$  satisfies  $\mathbf{V} = \mathbf{L}$  is much longer than I would like because the metamathematical techniques that are used to abbreviate textbook proofs are not available.

Once we have completed the first three tasks, we should be able to conclude that if ZF is consistent, then so is ZFL. (Also, from the fourth task, if ZF is consistent, then so is ZFC.) This inference requires reasoning in the metatheory, which is not possible using Isabelle/ZF, so the machine formalization omits it. Standard treatments also gloss over this step, regarding it as obvious. Section 2.4 below expands on this issue.

### 2.1. The problem with class models

Because  $\mathbf{L}$  is a proper class, we cannot adopt the usual notion of satisfaction. To formalize the standard Tarski definition of truth [10, p. 60] requires first defining, in set theory, a set  $F$  to represent the syntax of first-order formulae.  $F$  is easily defined, either using Gödel-numbering or as a recursive data structure. If: (i)  $M$  is a set, (ii)  $p \in F$  represents a formula with  $k$  free variables, and (iii)  $m_1, \dots, m_k \in M$ , then  $M \models p(m_1, \dots, m_k)$  can be defined by recursion on the structure of  $p$ . If  $\mathbf{M}$  is a proper class, then the obvious definition of  $\mathbf{M} \models p(\vec{m})$  cannot be formalized in set theory; the environments that hold the bindings of free variables would have to belong to a function space whose range was all of  $\mathbf{M}$ .

Tarski's theorem on non-definability of truth [9, p. 41] asserts that no formula  $\chi(p)$  expresses  $\mathbf{V} \models p$ . If, for each formula  $\phi$ , we write  $\ulcorner \phi \urcorner$  for the corresponding element of  $F$ ,

then  $\psi \leftrightarrow \neg\chi(\ulcorner\psi\urcorner)$  is a theorem for some sentence  $\psi$ . Satisfaction cannot be defined, at least if  $\mathbf{M} = \mathbf{V}$ .

## 2.2. Relativization

Gödel instead expressed satisfaction for class models syntactically. This approach abandons the set  $F$  of formula representatives in favour of real formulae. Set theory uses a first-order language with no constant symbols, no function symbols, and no relation symbols other than  $\in$  and  $=$ . Variables are the only terms.

Gödel's key concept is *relativization*; see [5, p. 76] or, for a modern treatment, [9, p. 112]. If  $\mathbf{M}$  is a class and  $\phi$  is a formula, define  $\phi^{\mathbf{M}}$  recursively as follows:

$$\begin{aligned}(x = y)^{\mathbf{M}} & \text{ abbreviates } x = y; \\ (x \in y)^{\mathbf{M}} & \text{ abbreviates } x \in y; \\ (\phi \wedge \psi)^{\mathbf{M}} & \text{ abbreviates } \phi^{\mathbf{M}} \wedge \psi^{\mathbf{M}}; \\ (\neg\phi)^{\mathbf{M}} & \text{ abbreviates } \neg(\phi^{\mathbf{M}}); \\ (\exists x. \phi)^{\mathbf{M}} & \text{ abbreviates } \exists x. x \in \mathbf{M} \wedge \phi^{\mathbf{M}}.\end{aligned}$$

Dually,  $(\forall x. \phi)^{\mathbf{M}}$  abbreviates  $\forall x. x \in \mathbf{M} \rightarrow \phi^{\mathbf{M}}$ , if universal quantifiers are defined as usual. (When working in ZF, we should write  $\mathbf{M}(x)$  instead of  $x \in \mathbf{M}$  above.) Relativization bounds all quantifiers in  $\phi$  by  $\mathbf{M}$ . It is intuitively clear that  $\phi^{\mathbf{M}}$  expresses that  $\phi$  is true in  $\mathbf{M}$ . However, while the satisfaction relation ( $\models$ ) can be defined within set theory, relativization can only be defined in the metalanguage: it combines two arguments,  $\phi$  and  $\mathbf{M}$ , which lie outside ZF.

## 2.3. The formal treatment of terms

Despite the lack of terms in their formal language, set theorists use elaborate notational conventions. In other branches of mathematics, an expression like  $f(x)g(y) - h(x, y)$  means what it says: functions  $f$ ,  $g$  and  $h$  are applied, and the results are combined by multiplication and subtraction. In set theory, however, each expression  $E(x)$  abbreviates a formula  $\phi(x, y)$ , which reduces the meaning of  $y = E(x)$  to a combination of  $\in$  and  $=$ . For example, we can express the meaning of  $Y = A \cup B$  by the predicate *union* ( $A, B, Y$ ), defined by

$$\forall z. z \in Y \leftrightarrow z \in A \vee z \in B.$$

We can similarly define *inter* ( $A, B, Y$ ) to express  $Y = A \cap B$ . Combining these predicates gives meaning to more complex terms; for example,  $Y = (A \cup B) \cap C$  abbreviates to

$$\exists X. \text{union}(A, B, X) \wedge \text{inter}(X, C, Y).$$

Variable binding notation, ubiquitous in set theory, causes further complications. In  $\bigcup_{x \in A} B(x)$ , what is  $B$ ? Syntactically,  $B(x)$  is a term with parameter  $x$ , so we can take it as an abbreviation for some formula  $\phi(x, y)$ . But then  $\bigcup$  becomes an operation on formulae, rather than one on sets. An equally legitimate alternative [7, p. 34] is to regard  $B$  as a function in set theory – formally, the set of pairs  $\{\langle x, B(x) \rangle \mid x \in A\}$ .

Set theorists generally say little about these notational conventions, and act as if terms were meaningful in themselves. Relativization, however, forces us to make the translation from terms to formulae explicit. In the Isabelle formalization, I have defined relational

equivalents of dozens of term-formers. I have included a class argument in each one to perform relativization at the same time; we can express the relativized term  $((A \cup B) \cap C)^{\mathbf{M}}$  as follows.

$$\exists X \in \mathbf{M}. \text{union}(\mathbf{M}, A, B, X) \wedge \text{inter}(\mathbf{M}, X, C, Y)$$

The hardest tasks were: (i) to define relational equivalents of the complicated expressions generated by Isabelle/ZF for recursively defined sets and functions, and (ii) to cope with the sheer bulk of the definitions.

## 2.4. Gödel's claim viewed proof-theoretically

The purpose of relativization is to express claims of the form ‘ $\phi$  is true in  $\mathbf{M}$ ’. To prove that  $\mathbf{L}$  satisfies the ZF axioms and that  $\mathbf{V} = \mathbf{L}$ , we must prove that  $\phi^{\mathbf{L}}$  for each ZF axiom  $\phi$ , and we must prove  $(\mathbf{V} = \mathbf{L})^{\mathbf{L}}$ . Now we can consider Gödel's claim that from a contradiction in ZFL, a contradiction in ZF ‘could actually be constructed’. His claim is proof-theoretic. A contradiction in ZFL is a proof,  $\Pi$ , of  $\perp$  from finitely many ZF axioms and  $\mathbf{V} = \mathbf{L}$ .

$$\begin{array}{c} \phi_1 \quad \dots \quad \phi_n \quad \mathbf{V} = \mathbf{L} \\ \Pi \\ \perp \end{array}$$

Once we have proved that  $\mathbf{L}$  satisfies the axioms of ZFL, we have the  $n + 1$  proofs

$$\text{ZF} \vdash \phi_1^{\mathbf{L}} \quad \dots \quad \text{ZF} \vdash \phi_n^{\mathbf{L}} \quad (\mathbf{V} = \mathbf{L})^{\mathbf{L}}.$$

Verifying Gödel's claim reduces to showing that we can always construct a proof  $\Pi^{\mathbf{L}}$  of  $\perp^{\mathbf{L}}$  from the following relativized premises.

$$\begin{array}{c} \phi_1^{\mathbf{L}} \quad \dots \quad \phi_n^{\mathbf{L}} \quad (\mathbf{V} = \mathbf{L})^{\mathbf{L}} \\ \Pi^{\mathbf{L}} \\ \perp^{\mathbf{L}} \end{array}$$

Then we get a proof of  $\text{ZF} \vdash \perp^{\mathbf{L}}$ , which is just  $\text{ZF} \vdash \perp$ .

So how *do* we obtain  $\Pi^{\mathbf{L}}$  from  $\Pi$ ? To be concrete, suppose that we are working with a natural deduction formalization of first-order logic. By the normal form theorem [20], since the conclusion of the proof is atomic, we can assume that  $\Pi$  applies only elimination rules. We must modify  $\Pi$  so that it accepts relativized versions of its premises and delivers a relativized version of its conclusion. The only hard cases involve quantifiers. Where  $\Pi$  applies the existential elimination rule to  $\exists x. \phi(x)$ , it delivers the formula  $\phi(x)$  to the rest of the proof. (Assume that  $x$  has already been renamed, if necessary.) At the corresponding position,  $\Pi^{\mathbf{L}}$  should apply the existential and conjunction elimination rules to  $\exists x. x \in \mathbf{L} \wedge \phi(x)$ , delivering the formulae  $x \in \mathbf{L}$  and  $\phi(x)$  to the rest of the proof.

Universal quantifiers require a bit more work. First, recall that the language of set theory has no terms other than variables. Where  $\Pi$  applies the universal elimination rule to  $\forall x. \phi(x)$ , it delivers the formula  $\phi(y)$  to the rest of the proof, where  $y$  is a variable. At the corresponding position,  $\Pi^{\mathbf{L}}$  should apply the existential and conjunction elimination rules to  $\forall x. x \in \mathbf{L} \rightarrow \phi(x)$ . But before it can deliver the formula  $\phi(y)$ , it requires a proof that  $y \in \mathbf{L}$ . We will indeed have  $y \in \mathbf{L}$  if the variable  $y$  is obtained by a previous existential elimination, but what if  $y$  was chosen arbitrarily? We can handle such cases by inserting at this point an application of the empty set axiom, which will yield a new variable (say  $z$ ) and the assumption that  $z \in \mathbf{L}$ . Intuitively, we are replacing all free variables in  $\Pi$  by 0.

The sketchy argument above cannot be called a rigorous proof of Gödel's claim, but it is more detailed than standard expositions of Gödel's proof. Kunen relegates the relevant lemma to an appendix, and for the proof he merely remarks 'Similar to the easy direction of the Gödel completeness theorem' [9, p. 141]. To Gödel, it was all presumably trivial. I have not formalized the argument in Isabelle/ZF because that would require formalizing the metatheory.

## 2.5. Defining the class $\mathbf{L}$

The equation  $\mathbf{V} = \bigcup_{\alpha \in \mathbf{ON}} V_\alpha$  expresses the universe of sets as the union of the cumulative hierarchy  $\{V_\alpha\}_{\alpha \in \mathbf{ON}}$ , which is recursively defined by  $V_0 = 0$ ,  $V_{\alpha+1} = \mathcal{P}(V_\alpha)$  and  $V_\alpha = \bigcup_{\xi < \alpha} V_\xi$  when  $\alpha$  is a limit ordinal. We obtain  $\mathbf{L}$  by a similar construction, replacing the powerset operator  $\mathcal{P}$  by the *definable powerset* operator,  $\mathcal{D}$ . Essentially,  $\mathcal{D}(A)$  yields the set of all subsets of  $A$  that can be defined by a formula taking parameters over  $A$ . If we define the set  $F$  of formulae and the satisfaction relation  $A \models p$  as outlined above, then we can make the definition

$$\mathcal{D}(A) = \{X \in \mathcal{P}(A) \mid \exists p \in F. \exists a_1 \dots a_k \in A. X = \{x \in A \mid A \models p(x, a_1, \dots, a_k)\}\}.$$

(The ellipsis can be eliminated in favour of lists over  $A$ .) Finally, we define the *constructible universe*:  $\mathbf{L} = \bigcup_{\alpha \in \mathbf{ON}} L_\alpha$ , where  $L_0 = 0$ ,  $L_{\alpha+1} = \mathcal{D}(L_\alpha)$  and  $L_\alpha = \bigcup_{\xi < \alpha} L_\xi$  when  $\alpha$  is a limit.

Kunen proves that  $\mathbf{L}$  satisfies the ZF axioms, remarking: 'only the Comprehension Axiom required any work' [9, p. 170]. His remark applies to the Isabelle/ZF proofs.  $\mathbf{L}$  inherits most of the necessary properties from  $\mathbf{V}$ . Even the axiom scheme of replacement can be proved as the theorem *replacement* ( $L, P$ ); the proof is independent of the formula  $P$ . However, the proof of comprehension for the formula  $\phi$  requires an instance of the reflection theorem for  $\phi$ , which requires recursion over the structure of  $\phi$ . Each instance of comprehension therefore has a different proof from the ZF axioms. At the metalevel, of course, all these proofs are instances of a single algorithm. For Isabelle/ZF, this means that each instance of comprehension must be proved separately, although the proof scripts are nearly identical.

## 2.6. Absoluteness: proving $(\mathbf{V} = \mathbf{L})^{\mathbf{L}}$

Proving that  $\mathbf{L}$  satisfies  $\mathbf{V} = \mathbf{L}$  is a key part of the proof, and despite first appearances, it is far from trivial. It amounts to saying that the construction of  $\mathbf{L}$  is idempotent:  $\mathbf{L}^{\mathbf{L}} = \mathbf{L}$ . The underlying concept is called *absoluteness*, which expresses that a given operator or formula behaves the same in a class model  $\mathbf{M}$  as it does in  $\mathbf{V}$ , the universe. A class  $\mathbf{M}$  is *transitive* if  $x \in \mathbf{M}$  implies that  $x \subseteq \mathbf{M}$ , and we shall be concerned only with transitive models below.

Most constructions are absolute. The empty set can only be a set having no elements, and  $A \subseteq B$  can only mean that every element of  $A$  belongs to  $B$ . If  $A$  and  $B$  are sets, then their union can only be the set containing precisely the elements of those sets. Many complicated notions are also absolute: domains and ranges of relations, bijections, well-orderings, order-isomorphisms, and ordinals. With some effort, we can show the absoluteness of recursively defined data structures and functions.

Powersets, except in trivial cases, are not absolute. For example,  $\mathcal{P}(\omega)$  might contain subsets of the natural numbers that cannot be shown to exist. The function space  $A \rightarrow B$  is not absolute because of the obvious connection between  $\mathcal{P}(A)$  and  $A \rightarrow \{0, 1\}$ . More subtly, cardinality is not absolute: if  $M$  is a countable model of set theory, and  $\alpha$  is an

uncountable cardinal according to  $M$ , then obviously  $\alpha$  must really be countable, with the bijections between  $\alpha$  and  $\omega$  lying outside  $M$ . This situation is called Skolem's paradox [9, p. 141].

Metamathematical arguments are an efficient means of proving absoluteness. For example, any concept that is provably equivalent (in ZF) to a formula involving only bounded quantifiers is absolute [9, p. 119]. This is the class of  $\Delta_0^{\text{ZF}}$  formulae. The larger class of  $\Delta_1^{\text{ZF}}$  formulae can also be shown to be absolute. Unfortunately, all such arguments are beyond our reach unless we formalize the metatheory.

## 2.7. The consequences of $\mathbf{V} = \mathbf{L}$

Once we have proved that  $\mathbf{L}$  is absolute, we obtain  $\text{ZF} \vdash (\mathbf{V} = \mathbf{L})^{\mathbf{L}}$ . We can then investigate the consequences of assuming that  $\mathbf{V} = \mathbf{L}$ . To prove the axiom of choice, it suffices to prove that every set can be well-ordered. The key step, given a well-ordering of  $A$ , is to construct a well-ordering of  $\mathcal{D}(A)$ . It comes from the lexicographic ordering on tuples  $\langle p, a_1, \dots, a_k \rangle$  for  $p \in F$  and  $a_1, \dots, a_k \in A$ . So if  $L_\alpha$  is well-ordered, so is  $L_{\alpha+1}$ . By transfinite induction, each level of the construction of  $\mathbf{L}$  is well-ordered.

The axiom  $\mathbf{V} = \mathbf{L}$  is very strong. Gödel proved that it implies the generalized continuum hypothesis. Jensen later proved that it implies the combinatorial principle known as  $\diamond$ , and it has many additional consequences. It is important to note, however, that such proofs are entirely separate from that of  $\text{ZF} \vdash (\mathbf{V} = \mathbf{L})^{\mathbf{L}}$ . We prove  $\text{ZFL} \vdash \text{AC}$ ,  $\text{ZFL} \vdash \text{GCH}$  and  $\text{ZFL} \vdash \diamond$ , but we do not prove  $\text{ZF} \vdash \text{AC}^{\mathbf{L}}$ ,  $\text{ZF} \vdash \text{GCH}^{\mathbf{L}}$  and  $\text{ZF} \vdash \diamond^{\mathbf{L}}$ . Those results, if we want them, are most easily obtained in the metatheory, using the general fact that if  $\phi \vdash \psi$ , then  $\phi^{\mathbf{L}} \vdash \psi^{\mathbf{L}}$ .

## 3. Introduction to the Isabelle/ZF formalization

Isabelle [11, 14] is an interactive theorem prover that supports a variety of logics, including set theory and higher-order logic. Isabelle provides automatic tools for simplification and logical reasoning. They can be combined with single-step inferences using a traditional tactical style or as structured proof texts. The **Proof General** user interface provides an effective interactive environment. Isabelle has been applied to a huge number of verification tasks, including the semantics of the Java language [21] and the correctness of cryptographic protocols [16]. Most of these proofs use Isabelle/HOL, the version of Isabelle for higher-order logic. Isabelle/HOL's polymorphic type system is ideal for modelling problems in computer science.

Isabelle also supports Zermelo–Fraenkel set theory. Formalized material includes the traditional concepts of functions, ordinals, order types and cardinals. Isabelle/ZF also accepts definitions of recursive functions and data structures; in this it resembles other computational logics, with the important difference of being typeless. Some problems do call for a typeless logic. Isabelle/ZF is also good for investigating foundational issues, and, of course, for formalizing proofs in axiomatic set theory.

Previously published work on Isabelle/ZF describes its basic development [13] and its treatment of recursive functions [13] and inductive definitions [17]. Another paper describes proofs drawn from set theory textbooks [19]. Particularly noteworthy are the proofs of equivalence between various formulations of the axiom of choice. Those proofs, formalized by Grąbaczewski, are highly technical, demonstrating that advanced set theory proofs can be replicated in Isabelle/ZF, given enough time and effort. That is precisely why we should



investigate Gödel's proof of the relative consistency of AC: much of the reasoning takes place outside set theory.

The previous section has presented many reasons why we should formalize Gödel's proof directly in the metatheory. That strategy does not require a set-theory prover. We could use any system that lets us define the first-order formulae, the set-theory axioms, and the set of theorems derivable from any given axioms. We would enjoy a number of advantages.

- Relativization could easily be defined by recursion on the structure of formulae.
- Metatheorems about absoluteness (for example, that all  $\Delta_0^{\text{ZF}}$  formulae are absolute) could be proved and used to obtain simple proofs of many absoluteness results.
- The constructiveness of the consistency result could be stated and proved.

However, the metatheoretical strategy also presents difficulties. We would have to work in the pure language of set theory, which reduces all concepts to membership and equality, and is unreadable; an alternative would be to formalize the familiar term language. We would constantly be reasoning about an explicitly formalized inference system for ZF, rather than using our prover's built-in reasoning tools. I believe that this strategy would involve as much work as the strategy I have adopted, although the work would be distributed differently.

The choice resembles the standard one that we face when we model a formal language: shall we adopt a deep or a shallow embedding? A shallow embedding maps phrases in the language to corresponding phrases in the prover's logic. It works well for reasoning about specific examples, but does not allow metareasoning (proofs about the language). A deep embedding involves formalizing the language's syntax and semantics in the prover's logic. The extra mechanism allows metareasoning, but complicates reasoning about specific examples. Compared with a shallow embedding, the strengths and weaknesses are exchanged.

I have chosen to formalize Gödel's theorem in set theory, minimizing any excursions into the metatheory. This strategy still requires the definition of relational equivalents for each element of set theory's term language, while limiting my exposure to unreadable relational formulae. After all, the critical proofs involve showing that various concepts are absolute, which means that they do not vary from one model of set theory to another. Each absoluteness proof justifies the replacement of some primitive of the relational language by its counterpart in the term language. Thus Isabelle's simplifier can transform relational formulae into ones using terms, exploiting the existing formalization of set theory.

This plan worked well for basic concepts such as union, intersection, relation, function, domain, range, image, inverse image and even ordinal. The absoluteness proofs for well-orderings, recursive functions and recursive data types were, however, harder.

- If a concept is defined in terms of non-absolute primitives, such as powerset, it must be proved equivalent to a suitable alternative definition.
- Much of the theory of well-founded recursion must be formalized from scratch in the relational language.
- Higher-order functions complicate the relational language.
- Recursive functions generate complicated fixedpoint definitions that must be manually converted into relational form.



#### 4. Relativization and absoluteness: the basics

The first step is to define the relational language, introducing predicates for all the basic concepts of set theory. Each predicate takes a class as an argument, so that it can express relativization. This relational language will later allow appeals to the reflection theorem. Space permits only a few of the predicates to appear below. Note that the class quantifications  $\forall x \in \mathbf{M}$  and  $\exists x \in \mathbf{M}$  are written  $\forall x[M]$  and  $\exists x[M]$  in Isabelle. For example,  $\forall x[M]. P(x)$  is definitionally equivalent to  $\forall x. M(x) \longrightarrow P(x)$ .

All Isabelle definitions in this paper are indicated by a vertical line, as shown below.

##### 4.1. From the empty set to functions

We begin with definitions of trivial concepts such as the empty set and the subset relation. A set  $z$  is empty if it has no elements.

$\text{"empty}(M, z) == \forall x[M]. x \notin z"$	
$\text{"subset}(M, A, B) == \forall x[M]. x \in A \longrightarrow x \in B"$	

A set  $z$  is the unordered pair of  $a$  and  $b$  if it contains those two sets and no others. The Kuratowski definition of ordered pairs  $\langle a, b \rangle = \{\{a, a\}, \{a, b\}\}$  is then expressed using the predicate  $\text{upair}$ .

$\text{"upair}(M, a, b, z) == a \in z \ \& \ b \in z \ \& \ (\forall x[M]. x \in z \longrightarrow x=a \mid x=b) "$	
$\text{"pair}(M, a, b, z) == \exists x[M]. \text{upair}(M, a, a, x) \ \& \ (\exists y[M]. \text{upair}(M, a, b, y) \ \& \ \text{upair}(M, x, y, z)) "$	

A set  $z$  is the union of  $a$  and  $b$  if it contains their elements and no others. The general union  $\bigcup(A)$ , also written as  $\bigcup\{X \mid X \in A\}$ , has an analogous definition.

$\text{"union}(M, a, b, z) == \forall x[M]. x \in z \longleftrightarrow x \in a \mid x \in b"$	
$\text{"big\_union}(M, A, z) == \forall x[M]. x \in z \longleftrightarrow (\exists y[M]. y \in A \ \& \ x \in y) "$	

A set  $z$  is the domain of the relation  $r$  if it consists of each element  $x$  such that  $\langle x, y \rangle \in r$  for some  $y$ .

$\text{"is\_domain}(M, r, z) == \forall x[M]. x \in z \longleftrightarrow (\exists w[M]. w \in r \ \& \ (\exists y[M]. \text{pair}(M, x, y, w))) "$	
---	--

##### 4.2. Relativizing the ordinals

Now we can define relational versions of ordinals and related concepts. The formalization is straightforward.

An ordinal is a transitive set of transitive sets.

$\text{"transitive\_set}(M, a) == \forall x[M]. x \in a \longrightarrow \text{subset}(M, x, a) "$	
$\text{"ordinal}(M, a) == \text{transitive\_set}(M, a) \ \& \ (\forall x[M]. x \in a \longrightarrow \text{transitive\_set}(M, x)) "$	

A limit ordinal is a non-empty, successor-closed ordinal.

$\text{"limit\_ordinal}(M, a) == \text{ordinal}(M, a) \ \& \ \sim \text{empty}(M, a) \ \& \ (\forall x[M]. x \in a \longrightarrow (\exists y[M]. y \in a \ \& \ \text{successor}(M, x, y))) "$	
---	--

A successor ordinal is any ordinal that is neither empty nor a limit.

```
"successor_ordinal(M, a) ==
  ordinal(M, a) & ~ empty(M, a) & ~ limit_ordinal(M, a) "
```

The set of natural numbers,  $\omega$ , is a limit ordinal that contains no limit ordinals.

```
"omega(M, a) ==
  limit_ordinal(M, a) & (∀ x[M]. x ∈ a → ~ limit_ordinal(M, x)) "
```

#### 4.3. Defining the Zermelo–Fraenkel axioms

Formally defining the ZF axioms relative to a class  $M$  lets us express the information that  $M$  satisfies those axioms. Each axiom is relativized so that all quantified variables range over  $M$ .

We begin with extensionality.

```
"extensionality(M) ==
  ∀ x[M]. ∀ y[M]. (∀ z[M]. z ∈ x ↔ z ∈ y) → x=y "
```

The separation axiom is also known as comprehension.

```
"separation(M, P) ==
  ∀ z[M]. ∃ y[M]. ∀ x[M]. x ∈ y ↔ x ∈ z & P(x) "
```

This yields a valid instance of separation only if the formula  $P$  obeys certain syntactic restrictions. All quantifiers in  $P$  must be relativized to  $M$ , and the free variables in  $P$  must range over elements of  $M$ . These restrictions prevent us from assuming separation as a scheme by leaving  $P$  as a free variable. We must separately note every instance of separation that we need. If it meets the syntactic restrictions, then later we shall be able to prove that  $\mathbf{L}$  satisfies it.

That looks bad when we recall that the native separation axiom in Isabelle/ZF, and the theorems using it, are schematic in  $P$ . However, if we formalize Bernays–Gödel set theory as a new Isabelle logic (creating the system Isabelle/BG), then the same problem occurs elsewhere. The analogue of separation in BG set theory is the general existence theorem, which is a metatheorem: proving each instance requires a separate construction. To compensate, at least BG has no axiom schemes.

The axioms of unordered pairs, unions and powersets all state that  $M$  is closed under the given operation, as follows.

```
"upair_ax(M) == ∀ x[M]. ∀ y[M]. ∃ z[M]. upair(M, x, y, z) "
"Union_ax(M) == ∀ x[M]. ∃ z[M]. big_union(M, x, z) "
"power_ax(M) == ∀ x[M]. ∃ z[M]. powerset(M, x, z) "
```

The foundation axiom states that every non-empty set has a  $\in$ -minimal element.

```
"foundation_ax(M) ==
  ∀ x[M]. (∃ y[M]. y ∈ x) → (∃ y[M]. y ∈ x & ~ (∃ z[M]. z ∈ x & z ∈ y)) "
```

We call a formula *univalent* over a set if it describes a class function on that set. The replacement axiom holds for univalent formulae as follows.

```
"univalent(M, A, P) ==
  ∀ x[M]. x ∈ A → (∀ y[M]. ∀ z[M]. P(x, y) & P(x, z) → y=z) "
"replacement(M, P) ==
  ∀ A[M]. univalent(M, A, P) →
    (∃ Y[M]. ∀ b[M]. (∃ x[M]. x ∈ A & P(x, b)) → b ∈ Y) "
```

Intuitively, if  $\mathcal{F}$  is a class function and  $A$  is a set, then replacement says that  $\mathcal{F}^{\ast}A$  (the image of  $A$  under  $\mathcal{F}$ ) is a set. However, the axiom formalized above is weaker: it merely asserts (relative to the class  $\mathbf{M}$ ) that  $\mathcal{F}^{\ast}A \subseteq B$  for some set  $B$ . To get the set that we really want, namely  $\mathcal{F}^{\ast}A$ , we must apply the axiom of separation to  $B$ . The weak form of replacement can be proved schematically for  $\mathbf{L}$ . The strong form cannot be proved schematically, because of its reliance on separation.

```
"strong_replacement(M,P) ==
  ∀ A[M]. univalent(M,A,P) →
  (∃ Y[M]. ∀ b[M]. b ∈ Y ↔ (∃ x[M]. x ∈ A & P(x,b))) "
```

#### 4.4. Introducing a transitive class model

The absoluteness proofs are carried out with respect to an arbitrary class model  $\mathbf{M}$ , although they are only needed for  $\mathbf{L}$ . Generalizing the proofs over other models has two advantages: it separates the absoluteness proofs from reasoning about  $\mathbf{L}$ , and it allows the proofs to be used with other class models.

Isabelle's locale mechanism [8] makes the generalization possible. A locale packages the many properties required of  $\mathbf{M}$ , creating a context in which they are implicitly available. A proof within a locale may refer to those properties, and also to other theorems proved in the same locale. A locale can extend an older one, creating a context that includes everything available in the ancestor locales.

```
locale M_trivial =
  fixes M
  assumes transM:      "[|y∈x; M(x)|] ⇒ M(y) "
  and upair_ax:        "upair_ax(M) "
  and Union_ax:        "Union_ax(M) "
  and power_ax:        "power_ax(M) "
  and replacement:    "replacement(M,P) "
  and M_nat [iff]:    "M(nat) "
```

The class  $\mathbf{M}$  is assumed to be transitive (*transM*), and to satisfy some relativized ZF axioms, such as unordered pairing (*upair\_ax*) and replacement. It contains the set of natural numbers, *nat* (which is also the ordinal  $\omega$ ). This locale does not assume any instances of separation.

#### 4.5. Easy absoluteness proofs

Here is a canonical example of an absoluteness result. The phrase **in** *M\_trivial* includes the lemma in the locale.

```
lemma (in M_trivial) empty_abs [simp]:
  "M(z) ⇒ empty(M,z) ↔ z=0 "
apply (simp add: empty_def)
apply (blast intro: transM)
done
```

The proof refers to the definition of empty set (*empty\_def*), and to the transitivity of  $M$  (the locale assumption *transM*); it uses *blast*, an automatic prover. The attribute [simp] declares *empty\_abs* as a simplification rule: the simplifier will replace any occurrence of *empty(M,z)* by *z=0*, provided that it can prove *M(z)*. From now on, usually just the statements of theorems will be shown, not header lines and proofs.

Here are some similar absoluteness results, also proved in locale *M\_trivial* and declared to the simplifier. Most have trivial proofs like the one shown above.

$$\begin{aligned}
 & "M(A) \implies \text{subset}(M, A, B) \longleftrightarrow A \subseteq B" \\
 & "M(z) \implies \text{upair}(M, a, b, z) \longleftrightarrow z = \{a, b\}" \\
 & "M(z) \implies \text{pair}(M, a, b, z) \longleftrightarrow z = \{a, b\}" \\
 & "[M(r); M(A); M(z)] \implies \text{image}(M, r, A, z) \longleftrightarrow z = r^{-1}A" \\
 & "[M(A); M(B); M(z)] \implies \text{cartprod}(M, A, B, z) \longleftrightarrow z = A \times B" \\
 & "[M(a); M(b); M(z)] \implies \text{union}(M, a, b, z) \longleftrightarrow z = a \cup b" \\
 & "[M(A); M(z)] \implies \text{big\_union}(M, A, z) \longleftrightarrow z = \text{Union}(A)"
 \end{aligned}$$

These theorems express absoluteness because the class  $M$  disappears from the right-hand side: the meanings of terms such as subset, image, and so on, are the same as their meanings in  $V$ . Each theorem also expresses the correctness of an element of the relational language, for example that *big\_union* captures the meaning of *Union*.

Absoluteness results involving ordinals are also easily proved, as follows.

$$\begin{aligned}
 & "M(a) \implies \text{ordinal}(M, a) \longleftrightarrow \text{Ord}(a)" \\
 & "M(a) \implies \text{limit\_ordinal}(M, a) \longleftrightarrow \text{Limit}(a)" \\
 & "M(a) \implies \text{successor\_ordinal}(M, a) \longleftrightarrow \\
 & \quad \text{Ord}(a) \ \& \ (\exists b[M]. \ a = \text{succ}(b))"
 \end{aligned}$$

Thus we see that the simplifier can rewrite relational formulae into term notation, provided that we are able to prove that they refer to elements of  $M$ . For this purpose, there are many results showing that  $M$  is closed under the usual set-theoretic constructions. In particular, we can use the separation axiom for a specific formula  $P$ , as shown here.

$$\begin{aligned}
 & "M(A) \implies M(\text{Union}(A))" \\
 & "[M(A); M(B)] \implies M(A \cup B)" \\
 & "[\text{separation}(M, P); M(A)] \implies M(\{x \in A. \ P(x)\})"
 \end{aligned}$$

Also useful are logical equivalences to simplify assertions involving  $M$ .

$$\begin{aligned}
 & "M(\{a, b\}) \longleftrightarrow M(a) \ \& \ M(b)" \\
 & "M(\langle a, b \rangle) \longleftrightarrow M(a) \ \& \ M(b)"
 \end{aligned}$$

#### 4.6. Absoluteness proofs assuming instances of separation

All the theorems shown above are proved without recourse to the axiom of separation. Obviously many set-theoretic operators are defined using separation – possibly in the guise of strong replacement – so we now extend locale *M\_trivial* accordingly.

```

locale M_basic = M_trivial +
assumes Inter_separation:
    "M(A)  $\implies$  separation(M,  $\lambda x. \ \forall y[M]. \ y \in A \longrightarrow x \in y$ )"
and Diff_separation:
    "M(B)  $\implies$  separation(M,  $\lambda x. \ x \notin B$ )"
and cartprod_separation:
    "[M(A); M(B)]
      $\implies$  separation(M,  $\lambda z. \ \exists x[M]. \ x \in A \ \& \$ 
        $(\exists y[M]. \ y \in B \ \& \ \text{pair}(M, x, y, z))$ )"
and image_separation:
    "[M(A); M(r)]
      $\implies$  separation(M,  $\lambda y. \ \exists p[M]. \ p \in r \ \& \$ 
        $(\exists x[M]. \ x \in A \ \& \ \text{pair}(M, x, y, p))$ )"
and converse_separation:
    "M(r)  $\implies$  separation(M,  $\lambda z. \ \exists p[M]. \ p \in r \ \& \$ 
        $(\exists x[M]. \ \exists y[M]. \ \text{pair}(M, x, y, p) \ \& \$ 
        $\text{pair}(M, y, x, z))$ )"
    
```

Only a few of the eleven instances of separation appear above. Omitted are the more complicated ones, for example those concerned with well-founded recursion.

By *Inter\_separation* it follows that  $M$  is closed under intersections.

**lemma** (in  $M\_basic$ ) *Inter\_closed*:  
 $"M(A) \implies M(Inter(A))"$

From the lemma declaration, you can see that the proof takes place in locale  $M\_basic$ . All results proved in locale  $M\_trivial$  remain available.

By *cartprod\_separation*, it follows that the class  $M$  is closed under Cartesian products. The proof is complicated because the powerset operator (which is not absolute) occurs in the definition. A trivial corollary is that  $M$  is closed under disjoint sums.

$"[M(A); M(B)] \implies M(A \times B)"$   
 $"[M(A); M(B)] \implies M(A + B)"$

I have devoted some effort to minimizing the number of instances of separation required. For example, the inverse image operator is expressed in terms of the image and converse operators. Then the domain and range operators can be expressed in terms of inverse image and image. We obtain five closure theorems from the two assumptions *image\_separation* and *converse\_separation*, as follows.

$"[M(A); M(r)] \implies M(r^{-1}A)"$   
 $"[M(A); M(r)] \implies M(r^{-}A)"$   
 $"M(r) \implies M(converse(r))"$   
 $"M(r) \implies M(domain(r))"$   
 $"M(r) \implies M(range(r))"$

These five operators are also absolute. Here is the result for *domain*.

$"[M(r); M(z)] \implies is\_domain(M, r, z) \longleftrightarrow z = domain(r)"$

Although we assume that  $M$  satisfies the powerset axiom, we cannot hope to prove that  $M(A) \implies M(Pow(A))$ . The powerset of  $A$  relative to  $M$  is smaller than the true powerset, containing only those subsets of  $A$  that belong to  $M$ . Similarly, we cannot show that  $M$  contains all functions from  $A$  to  $B$ . However, it holds for a finite case, essentially the set  $B^n$  of  $n$ -tuples, as shown here.

$"[n \in nat; M(B)] \implies M(n \rightarrow B)"$

This lemma will be needed later to prove the absoluteness of transitive closure.

#### 4.7. Some remarks about functions

In set theory, a function is a single-valued relation, and thus is a set of ordered pairs. Operators such as powerset and union, which apply to all sets, are not functions. (Strictly speaking, there are no operators in the formal language of set theory, since the only terms are variables.) Isabelle/ZF distinguishes functions from operators syntactically.

- The application of the function  $f$  to the argument  $x$  is written  $f \cdot x$ . On the other hand, application of an operator to its operand is written using parentheses, as in  $Pow(X)$ , or by using infix notation.

- Function abstraction over a set  $A$  is indicated by  $\lambda x \in A$ , and yields a set of pairs. For instance,  $\lambda x \in A. x$  denotes the identity function on  $A$ . Operators are essentially abstractions over the universe, as in  $\lambda x. Pow(Pow(x))$ . Abstraction can also express predicates; for instance,  $\lambda x. P(x) \ \& \ Q(x)$  is the conjunction of the two predicates  $P$  and  $Q$ .

Kunen [9, p. 14] defines function application in the usual way:  $f \cdot x$  is ‘the unique  $y$  such that  $\langle x, y \rangle \in f$ ’. Isabelle/ZF originally adopted a formal version of this definition, using a description operator [13, Section 7.5]. The relational version of the operator, namely  $\text{fun\_apply}(M, f, x, y)$ , held if the pair  $\langle x, y \rangle$  belongs to  $f$  for that unique  $y$ .

My original definitions of function application, in its infix and relational forms, both followed Kunen’s definition. However, the absoluteness theorem relating them was conditional on the function application’s being well-defined. That made it harder to simplify  $\text{fun\_apply}(M, f, x, y)$  to  $f \cdot x = y$ , and often forced proofs to include what was essentially type information.

Redefining function application by  $f \cdot x = \bigcup(f \cdot \{a\})$  solved these problems by eliminating the definite description. The new definition looks peculiar, but it agrees with the old one when the latter is defined. Its relational version is straightforward.

```
"fun_apply(M, f, x, y) ==
  (∃ xs[M]. ∃ fxs[M].
    upair(M, x, x, xs) & image(M, f, xs, fxs) & big_union(M, fxs, y)) "
```

Thus it follows that  $M$  is closed under function application, which is also absolute.

```
"[M(f); M(a)] ⟹ M(f ` a) "
"[M(f); M(x); M(y)] ⟹ fun_apply(M, f, x, y) ⟷ f ` x = y"
```

## 5. Well-founded recursion

The hardest absoluteness proofs concern recursion. Well-founded recursion is the most general form of recursive function definition. The proof that well-founded relations are absolute consists of several steps. Well-orderings, which are well-founded linear orderings, are somewhat easier to prove absolute.

### 5.1. Absoluteness of well-orderings

The concept of well-ordering is the first that we encounter of those whose absoluteness proof is hard. One direction is easy: if relation  $r$  well-orders  $A$ , then it also well-orders  $A$  relative to  $\mathbf{M}$ . For if every nonempty subset of  $A$  has an  $r$ -minimal element, then trivially so does every nonempty subset of  $A$  that belongs to  $\mathbf{M}$ ; this is Kunen’s [9, Lemma IV 3.14, p. 123]. For proving the converse direction, Kunen reasons in [9, Theorem IV 5.4, p. 127] that ‘every well-ordering is isomorphic to an ordinal’. We can obtain this result by showing that order types exist in  $\mathbf{M}$  and are absolute. The proof requires some instances of separation and replacement for  $\mathbf{M}$ .

The theory defines various properties of relations, relative to a class  $M$ . Transitivity, linearity, and other simple properties have the obvious definitions and are easily demonstrated to be absolute. The definition of ‘well-founded’ refers to the existence of  $r$ -minimal elements, as discussed above.

```
"wellfounded_on(M, A, r) ==
  ∀ x[M]. x ≠ 0 ⟶ x ⊆ A
  ⟶ (∃ y[M]. y ∈ x & ~ (∃ z[M]. z ∈ x & ⟨z, y⟩ ∈ r)) "
```

A well-ordering is a well-founded relation that is also linear and transitive.

```
"wellordered(M, A, r) ==
  transitive_rel(M, A, r) & linear_rel(M, A, r) &
  wellfounded_on(M, A, r) "
```

Kunen's [9, Lemma IV 3.14] takes the following form.

$$\text{"well\_ord}(A, r) \implies \text{wellordered}(M, A, r) \text{"}$$

The definition of order types is standard; see [9, Theorem I 7.6, p. 17]. We use replacement to construct a function that maps elements of  $A$  to ordinals, proving that its domain is the whole of  $A$ , and that each element of its range is an ordinal. Its range is the desired order type. But the construction must be done relative to  $M$ . In particular, when we need well-founded induction on  $r$ , we must apply a relativized induction rule, as follows.

$$\begin{aligned} &\text{"}[a \in A; \text{wellfounded\_on}(M, A, r); M(A); \\ &\quad \text{separation}(M, \lambda x. x \in A \longrightarrow \neg P(x)); \\ &\quad \forall x \in A. M(x) \ \& \ (\forall y \in A. \langle y, x \rangle \in r \longrightarrow P(y)) \longrightarrow P(x)] \\ &\implies P(a) \text{"} \end{aligned}$$

One premise is an instance of the separation axiom involving the negation of the induction formula. Each time that we apply induction, we must assume another instance of separation.

After about 250 lines of proof script, we arrive at [9, Theorem IV 5.4]. The notion of well-ordering is absolute.

$$\text{"}[M(A); M(r)] \implies \text{wellordered}(M, A, r) \longleftrightarrow \text{well\_ord}(A, r) \text{"}$$

Order types are absolute. That is, if  $f$  is an order-isomorphism from between  $(A, r)$  and some ordinal  $i$ , then  $i$  is the order type of  $(A, r)$ .

$$\begin{aligned} &\text{"}[ \text{wellordered}(M, A, r); f \in \text{ord\_iso}(A, r, i, \text{Memrel}(i)); \\ &\quad M(A); M(r); M(f); M(i); \text{Ord}(i) ] \implies i = \text{ordertype}(A, r) \text{"} \end{aligned}$$

These results are not required in the sequel, but I found their proofs a useful preparation for tackling the more general problem of well-founded recursion.

## 5.2. Functions defined by well-founded recursion are absolute

It is essential to show that functions can be defined by well-founded recursion in  $M$ , and that such functions are absolute. This is Kunen's [9, Theorem IV 5.6, p. 129].

Let  $r$  be a well-founded relation. If  $f$  is recursively defined over  $r$ , then  $f(a)$  is derived from  $a$  and from various  $f(y)$ , where  $y$  ranges over the set of  $r$ -predecessors of  $a$ . This set is just  $r^{-1}\{a\}$ , the inverse image of  $\{a\}$  under  $r$  or, more explicitly,  $\{y \mid \langle y, a \rangle \in r\}$ . Writing the body of  $f$  as  $H(x, g)$ , with free variables  $x$  and  $g$ , we get the recursion equation

$$f(a) = H(a, f \upharpoonright (r^{-1}\{a\})). \quad (1)$$

Note that  $f \upharpoonright (r^{-1}\{a\})$  denotes the function obtained by restricting  $f$  to  $r$ -predecessors of  $a$ .

If  $r$  and  $H$  are given, then the existence of a suitable function  $f$  follows by well-founded induction over  $r$ , as I have described in previous work [15]. I have had to repeat some of these proofs relative to  $M$ . The theorems may assume only the relativized assumption  $\text{wellfounded}(M, r)$ , which for the moment is weaker than  $\text{wf}(r)$ . About 200 lines of proof script are necessary, but fortunately much of this material is based on earlier proofs. We reach the following key result concerning the existence of recursive functions.

$$\begin{aligned} &\text{"}[ \text{wellfounded}(M, r); \text{trans}(r); \\ &\quad \text{separation}(M, \lambda x. \sim (\exists f[M]. \text{is\_recfun}(r, x, H, f))); \\ &\quad \text{strong\_replacement}(M, \lambda x z. \\ &\quad \quad \exists y[M]. \exists g[M]. z = \langle x, y \rangle \ \& \ \text{is\_recfun}(r, x, H, g) \ \& \ y = H(x, g)); \\ &\quad M(r); M(a); \\ &\quad \forall x[M]. \forall g[M]. \text{function}(g) \longrightarrow M(H(x, g))] \\ &\implies \exists f[M]. \text{is\_recfun}(r, a, H, f) \text{"} \end{aligned}$$



The predicate  $is\_recfun(r, a, H, f)$  expresses the fact that  $f$  satisfies the recursion equation (1) for the given relation  $r$  and body  $H$  for all  $r$ -predecessors of  $a$ . So the theorem states that if  $r$  is well-founded and transitive, then there exists  $f$  in  $M$  satisfying the recursion equation below  $a$ . Obviously,  $r$  and  $a$  must belong to the class  $M$ , which moreover must be closed under  $H$ . Two additional premises list instances of separation and replacement, which depend upon  $r$  and  $H$ . Before we can assume such instances, we must express them relative to  $M$ . That in turn requires a relativized version of  $is\_recfun$ .

```
"M_is_recfun(M, MH, r, a, f) ==
  ∀z[M]. z ∈ f ↔
    (∃x[M]. ∃y[M]. ∃xa[M]. ∃sx[M]. ∃r_sx[M]. ∃f_r_sx[M].
      pair(M, x, y, z) & pair(M, x, a, xa) & upair(M, x, x, sx) &
      pre_image(M, r, sx, r_sx) & restriction(M, f, r_sx, f_r_sx) &
      xa ∈ r & MH(x, f_r_sx, y))"
```

This definition is the translation of equation (1) into relational language. (Observe how quickly this language becomes unreadable.) In particular, the binary operator  $H$  becomes the ternary relation  $MH$ . The argument  $H$  makes  $is\_recfun$  a higher-order function, which complicates subsequent work. We cannot relativize  $is\_recfun$  once and for all, but if  $MH$  is expressed in relational language, then so is  $M\_is\_recfun$ .

The predicate  $relation2$  expresses that  $is\_f$  is the relational form of  $f$  over  $M$ .

```
"relation2(M, is_f, f) ==
  ∀x[M]. ∀y[M]. ∀z[M]. is_f(x, y, z) ↔ z = f(x, y)"
```

The predicate  $is\_wfrec$  expresses that  $z$  is computed from  $a$  and  $MH$  by well-founded recursion over  $r$ . The body of the definition expresses the existence of a function  $f$  satisfying equation (1), with  $z = H(a, f)$ .

```
"is_wfrec(M, MH, r, a, z) ==
  ∃f[M]. M_is_recfun(M, MH, r, a, f) & MH(a, f, z)"
```

We now reach two lemmas, stating that  $M\_is\_recfun$  and  $is\_wfrec$  behave as intended. The first result is the absoluteness of  $is\_recfun$ . Among the premises are that  $M$  is closed under  $H$ , and that  $MH$  is the relational form of  $H$ .

```
"[∀x[M]. ∀g[M]. function(g) → M(H(x, g)); M(r); M(a); M(f);
  relation2(M, MH, H)]
⇒ M_is_recfun(M, MH, r, a, f) ↔ is_recfun(r, a, H, f)"
```

Under identical premises, we get the following corollary.

```
"is_wfrec(M, MH, r, a, z) ↔ (∃g[M]. is_recfun(r, a, H, g) & z=H(a, g))"
```

### 5.3. Making well-founded recursion available

Mathematically speaking, we have already proved the absoluteness of well-founded recursion. Pragmatically speaking, unfortunately, more work must be done to package the results so that they can be used in formal proofs. In particular, we need a theorem relating the predicate  $is\_wfrec$ , defined above, with the function  $wfrec$  provided by Isabelle/ZF [15, Section 3.1];  $wfrec(r, a, H)$  denotes the result of  $f(a)$ , where  $f$  is the function with body  $H$  defined by recursion over  $r$ .

The development of well-founded recursion assumes  $r$  to be transitive. To apply well-founded recursion to other relations requires a theory of transitive closure. Isabelle/ZF

defines the transitive closure of a relation inductively [15, Section 2.5]. Inductive definitions are abstract and elegant, but they do not lend themselves to absoluteness proofs because they use the powerset operator. We must find an alternative definition, and an obvious one is based on the intuition that

$$x <^* y \iff x = s_0 < s_1 < \dots < s_n = y.$$

The sequence  $s_0, s_1, \dots, s_n$  can be modelled as a finite function: as noted in Section 4.6, finite functions are absolute. From  $x <^* y$  it is trivial to define the transitive closure,  $x <^+ y$ . In the definition below,  $f$  is the sequence and  $A$  is intended to represent the field of  $r$ .

```
"rtrancl_alt(A,r) ==
  {p ∈ A*A. ∃ n ∈ nat. ∃ f ∈ succ(n) -> A.
    (∃ x y. p = ⟨x,y⟩ & f'0 = x & f'n = y) &
    (∀ i ∈ n. ⟨f'i, f'succ(i)⟩ ∈ r)}"
```

It is easy to prove that this definition coincides with Isabelle/ZF's inductive one, which reads as follows.

```
"rtrancl_alt(field(r),r) = r^*"
```

Since every concept used in the new definition is absolute, we merely have to relativize this definition to  $M$ , defining  $rtran\_closure\_mem(M,A,r,p)$  to hold when  $p$  is an element of  $rtrancl\_alt(A,r)$ . I omit the definition because the relational language is unreadable. We cannot even use the constant  $0$ , but must introduce a variable  $zero$  and constrain it by  $empty(M,zero)$ .

The next two predicates relativize the reflexive-transitive and transitive closure of a relation.

```
"rtran_closure(M,r,s) ==
  ∀ A[M]. is_field(M,r,A) →
    (∀ p[M]. p ∈ s ↔ rtran_closure_mem(M,A,r,p))"
"tran_closure(M,r,t) ==
  ∃ s[M]. rtran_closure(M,r,s) & composition(M,r,s,t)"
```

Once we assume an instance of separation involving  $rtran\_closure\_mem$ , the closure and absoluteness results follow directly.

```
"M(r) ⇒ M(rtrancl(r))"
"⟦M(r); M(z)⟧ ⇒ rtran_closure(M,r,z) ↔ z = rtrancl(r)"
"M(r) ⇒ M(trancl(r))"
"⟦M(r); M(z)⟧ ⇒ tran_closure(M,r,z) ↔ z = trancl(r)"
```

If a relation is well-founded, then so is its transitive closure. The following lemma is useful because at this point we do not know that  $wellfounded(M,r)$  is equivalent to  $wf(M,r)$ .

```
"⟦wellfounded(M,r); M(r)⟧ ⇒ wellfounded(M,r^+)"
```

After about 130 lines of proof script, we arrive at some important theorems. One asserts absoluteness, relating the predicate  $is\_wfrec$  with the operator  $wfrec$ .

```
"⟦wf(r); trans(r); relation(r); M(r); M(a); M(z);
  wfrec_replacement(M,MH,r); relation2(M,MH,H);
  ∀ x[M]. ∀ g[M]. function(g) → M(H(x,g))⟧
⇒ is_wfrec(M,MH,r,a,z) ↔ z = wfrec(r,a,H)"
```

Another states that the class  $M$  is closed under well-founded recursion.

```
"[[wf(r); trans(r); relation(r); M(r); M(a);
  wfrec_replacement(M,MH,r); relation2(M,MH,H);
  ∀x[M]. ∀g[M]. function(g) → M(H(x,g))]]
⇒ M(wfrec(r,a,H))"
```

The theorems fortunately require identical instances of replacement. Both theorems assume that  $\text{trans}(r)$  holds; omitted here are more general theorems that relax the assumption of transitivity.

Both theorems use the predicate  $\text{wfrec\_replacement}$  to express a necessary instance of replacement. Its arguments are the ternary predicate  $MH$ , which represents the body of the recursive function, and the well-founded relation  $r$ .

```
"wfrec_replacement(M,MH,r) ==
  strong_replacement(M,
    λx z. ∃y[M]. pair(M,x,y,z) & is_wfrec(M,MH,r,x,y))"
```

## 6. Defining first-order formulae and the class **L**

We pause from proving absoluteness results in order to consider our main objective, namely the class **L** and its properties. The most logical order of presentation might have been to develop **L** first, and then to prove that constructibility is absolute. The order of presentation adopted here better represents how I actually carried out the proofs. Kunen similarly presents general absoluteness results before he introduces **L**.

### 6.1. Internalized first-order formulae

The idea of **L** is to introduce, at each stage, the sets that can be defined from existing ones by a first-order formula with parameters. Neither Gödel [5] nor Kunen actually uses first-order formulae, preferring more abstract constructions that achieve the goal more easily. However, Isabelle/ZF's recursive datatype package automates the task of defining the set of first-order formulae and the satisfaction relation on them. Gödel's earlier proof [6] also uses first-order formulae.

The obvious representation of first-order formulae is de Bruijn's [1], where there are no variable names. Instead, each variable reference is a non-negative integer, where zero refers to the innermost quantifier, and larger numbers refer to enclosing quantifiers. If the integer is greater than the number of enclosing quantifiers, then it is a free variable. This representation eliminates the danger of name confusion. It is particularly useful for formulae with parameters, since their order is determined numerically, rather than by name.

```
datatype "formula" = Member ("x ∈ nat", "y ∈ nat")
  | Equal ("x ∈ nat", "y ∈ nat")
  | Nand ("p ∈ formula", "q ∈ formula")
  | Forall ("p ∈ formula")
```

Having only four cases simplifies the relativization of functions on formulae. All propositional connectives are expressed in terms of  $\text{Nand}$ .

```
"Neg(p) == Nand(p,p) "
"And(p,q) == Neg(Nand(p,q)) "
"Or(p,q) == Nand(Neg(p),Neg(q)) "
"Implies(p,q) == Nand(p,Neg(q)) "
"Iff(p,q) == And(Implies(p,q), Implies(q,p)) "
"Exists(p) == Neg(Forall(Neg(p))) "
```

## 6.2. The satisfaction relation

Satisfaction is a primitive recursive function on formulae. Thanks to the nameless representation, interpretations are simply lists rather than functions from variable names to values. The familiar list function *nth*, defined below, looks up variables in interpretations.

```
"nth(0, Cons(a, l)) = a"
"nth(succ(n), Cons(a,l)) = nth(n,l)"
"nth(n, Nil) = 0"
```

The second of these equations is subject to the condition  $n \in \text{nat}$ . Note that element zero is the head of the list. Another useful function is *bool\_of\_o*, which converts a truth value to an integer, as follows.

```
"bool_of_o(P) == (if P then 1 else 0)"
```

This conversion is necessary because Isabelle/ZF is based on first-order logic. Formulae are not values, so we encode them using integers. We thus define a recursive predicate as a recursive integer-valued function. We are now able to define the function *satisfies*, which takes a set (the domain of discourse), a formula and an interpretation (written *env* for environment). It returns 1 or 0, depending upon whether or not the formula evaluates to true or false.

```
"satisfies(A, Member(x,y)) =
  (lenv ∈ list(A). bool_of_o (nth(x, env) ∈ nth(y, env)))"
"satisfies(A, Equal(x,y)) =
  (lenv ∈ list(A). bool_of_o (nth(x, env) = nth(y, env)))"
"satisfies(A, Nand(p,q)) =
  (lenv ∈ list(A). not ((satisfies(A,p)'env) and
    (satisfies(A,q)'env)))"
"satisfies(A, Forall(p)) =
  (lenv ∈ list(A). bool_of_o
    (∀x∈A. satisfies(A,p)'(Cons(x, env)) = 1))"
```

The abstraction and explicit function applications involving environments are necessary because the environments can vary in the recursive calls. The last line of *satisfies* deserves attention. The universal formula *Forall*(*p*) evaluates to 1 just if *p* evaluates to 1 in every environment obtainable from *env* by adding an element of *A*. Such environments have the form *Cons*(*x*, *env*) for  $x \in A$ .

The satisfaction predicate, *sats*, is a macro that refers to the function *satisfies*.

```
translations "sats(A,p,env)" == "satisfies(A,p)'env = 1"
```

The satisfaction predicate enjoys a number of properties that relate the internalized formulae to real formulae. All the equivalences are subject to the typing condition  $\text{env} \in \text{list}(A)$ . For example, the membership and equality relations behave as they should.

```
"sats(A, Member(x,y), env) ↔ nth(x, env) ∈ nth(y, env)"
"sats(A, Equal(x,y), env) ↔ nth(x, env) = nth(y, env)"
```

The propositional connectives also work.

```
"sats(A, Neg(p), env) ↔ ~ sats(A,p,env)"
"(sats(A, And(p,q), env)) ↔ sats(A,p,env) & sats(A,q,env)"
"(sats(A, Or(p,q), env)) ↔ sats(A,p,env) | sats(A,q,env)"
```

Quantifiers work too. Notice how the environment is extended.

```
"sats(A, Exists(p), env) ↔ (∃x∈A. sats(A, p, Cons(x, env)))"
```

### 6.3. The arity of a formula

The arity of a formula is, intuitively, its set of free variables. In  $\text{sats}(A, p, \text{env})$ , if the arity of  $p$  does not exceed the length of  $\text{env}$ , then the environment supplies values to all of  $p$ 's free variables.

Take each de Bruijn reference, adjusted for the depth of quantifier nesting at that point; the arity is the maximum of the resulting values. The recursive definition of function *arity* is simpler than this description.

```
"arity(Member(x,y)) = succ(x) ∪ succ(y) "
"arity(Equal(x,y)) = succ(x) ∪ succ(y) "
"arity(Nand(p,q)) = arity(p) ∪ arity(q) "
"arity(Forall(p)) = Arith.pred(arity(p)) "
```

Note that  $m \cup n = \max\{m, n\}$  in set theory, and that *Arith.pred* denotes the predecessor function. Trivial corollaries of this definition tell us how to compute the arities of other connectives as follows.

```
"arity(Neg(p)) = arity(p) "
"arity(And(p,q)) = arity(p) ∪ arity(q) "
```

The following result is more interesting. Extra items in the environment (exceeding the arity) are ignored. Here '@' is the list 'append' operator, so  $\text{env@extra}$  is  $\text{env}$  with additional items added.

```
"[arity(p) ≤ length(env); p ∈ formula;
  env ∈ list(A); extra ∈ list(A)]
⇒ sats(A, p, env@extra) ↔ sats(A, p, env) "
```

### 6.4. Renaming (renumbering) free variables

If  $A$  is a set, then the subset

$$\{x \in A \mid \phi(x, a_1, \dots, a_m)\}$$

is determined by the choice of  $\phi$  and of the parameters  $a_1, \dots, a_m$ , which are elements of  $A$ . These are the definable subsets of  $A$ .

Now, consider the problem of showing that the definable sets are closed under intersection. Suppose that another subset of  $A$  is defined by a formula  $\psi$  and parameters  $a_{m+1}, \dots, a_{m+n}$ .

$$\{x \in A \mid \psi(x, a_{m+1}, \dots, a_{m+n})\}$$

Then their intersection can presumably be defined by

$$\{x \in A \mid \phi(x, a_1, \dots, a_m) \wedge \psi(x, a_{m+1}, \dots, a_{m+n})\}.$$

Our aim is to regard the conjunction  $\phi \wedge \psi$  as having the free variables  $x, a_1, \dots, a_n$ . The occurrences of  $x$  in both formulae must be identified, while the parameter lists of the two formulae must be kept disjoint. To achieve our aim, we may need to rename one of the formula's free variables.

The de Bruijn representation refers to variables by number rather than by name. The variables shown as  $x$  above always have the de Bruijn index zero, so they will be identified automatically. We keep the parameter lists disjoint by renumbering the free variables in one of the formulae. Since  $x$  must be left alone, we renumber only the variables having an index greater than zero.

Renumbering functions are often necessary with the de Bruijn approach, though normally they rename variables during substitution. When efficiency matters, the renumbering functions take an argument specifying what number should be added to the variables. Here, the definitions are used for reasoning about, rather than for execution, so for us renaming means adding one; repeating this allows renaming by larger integers. In the following definitions,  $nq$  refers to the number of quantifiers enclosing the current point. Any de Bruijn index smaller than  $nq$  must not be renamed.

#### 6.4.1. The renaming function

First, we need a one-line function that renames a de Bruijn variable.

```
"incr_var(x,nq) == if x<nq then x else succ(x)"
```

Now we can define the main renaming function. As with *satisfies* above, abstraction and explicit function applications are necessary: the argument  $nq$  ('nesting of quantifiers') varies in the recursive calls. In the *Member* and *Equal* cases, the variables are simply renamed. The *Nand* case recursively renames the subformulae using the same nesting depth, while the *Forall* case renames its subformula using an increased nesting depth.

```
"incr_bv(Member(x,y)) =
  (λnq ∈ nat. Member (incr_var(x,nq), incr_var(y,nq)))"

"incr_bv(Equal(x,y)) =
  (λnq ∈ nat. Equal (incr_var(x,nq), incr_var(y,nq)))"

"incr_bv(Nand(p,q)) =
  (λnq ∈ nat. Nand (incr_bv(p)'nq, incr_bv(q)'nq))"

"incr_bv(Forall(p)) =
  (λnq ∈ nat. Forall (incr_bv(p) ' succ(nq)))"
```

Recall the example at the start of this section, concerning a set defined by the conjunction  $\phi \wedge \psi$ . If we are to conjoin the formulae  $\phi$  and  $\psi$  and combine their sets of parameters, then we need to ensure that some of the parameters are visible only to  $\phi$ , while the rest are visible only to  $\psi$ . The following lemma makes this possible.

```
"[p ∈ formula; bvs ∈ list(A); env ∈ list(A); x ∈ A]
⇒ sats(A, incr_bv(p) ' length(bvs), bvs @ Cons(x,env)) ↔
  sats(A, p, bvs@env)"
```

For the intuition, suppose that  $bvs$  is the list  $[x_0, \dots, x_{m-1}]$  (and therefore has length  $m$ ). Then the conclusion essentially says that

$$\text{sats}(A, \text{incr\_bv}(p) \text{ ' } m, [x_0, \dots, x_{m-1}, x, x_m, \dots, x_n]) \longleftrightarrow \text{sats}(A, p, [x_0, \dots, x_{m-1}, x_m, \dots, x_n])$$

and thus the renaming allows an additional value to be put into the environment at position  $m$ . The renamed formula will ignore the new value. By repeated renaming, we can construct a formula that will ignore a section of the parameter list that is intended for another formula.

The next result describes the obvious relationship between *arity* and renaming. Renaming increases a formula's arity by one, unless the variable being renamed does not exist, when renaming has no effect.

```
"[p ∈ formula; n ∈ nat]
⇒ arity(incr_bv(p) ' n) =
  (if n < arity(p) then succ(arity(p)) else arity(p))"
```

Considering how trivial the notion of arity is, many proofs about it (including this one) are complicated by innumerable case splits. Getting the simplifier to prove most of them automatically requires some ingenuity. Many other tiresome proofs about arities are omitted here.

#### 6.4.2. Renaming all but the first bound variable

One more thing is needed before we can define sets using conjunctions. As discussed at the beginning of Section 6.4, when a formula  $\phi$  defines a set, the variable with de Bruijn index zero gives the extension of that set, while the remaining free variables serve as parameters. Therefore, our basic renaming operator must only rename variables having a de Bruijn index of one or more.

```
| "incr_bv1(p) == incr_bv(p) '1"
```

Finally, we reach a lemma justifying our intended use of renaming.

```
"[[p ∈ formula; bvs ∈ list(A); x ∈ A; env ∈ list(A);
  length(bvs) = n]]
⇒ sats(A, iterates(incr_bv1, n, p), Cons(x, bvs@env)) ↔
   sats(A, p, Cons(x, env))"
```

If the environment has an initial segment  $bvs$  of length  $n$ , and if we apply the  $incr\_bv1$   $n$  times, then the modified formula ignores the  $bvs$  part. However, the renamed and original formulae agree on the first element of the environment, shown above as  $x$ .

#### 6.5. The definable powerset operation

The definable powerset operator is called  $DPow$ .

```
"DPow(A) == {X ∈ Pow(A).
  ∃ env ∈ list(A). ∃ p ∈ formula.
    arity(p) ≤ succ(length(env)) &
    X = {x ∈ A. sats(A, p, Cons(x, env))}}"
```

A set  $X$  belongs to  $DPow(A)$ , provided that there is an environment  $env$  (a list of values drawn from  $A$ ) and a formula  $p$ . The constraint  $arity(p) \leq succ(length(env))$  indicates that the environment should interpret all but one of  $p$ 's free variables. The variable whose de Bruijn index is zero determines the extension of  $x$  via the satisfaction relation:  $sats(A, p, Cons(x, env))$ . You may want to compare this with the informal discussion in the previous section, or with of Kunen's [9, Definition VI 1.1, p. 165].

Some consequences of this definition are easy to prove. The empty set is defined by the predicate  $\lambda x. x \neq x$ , and singleton sets by  $\lambda x. x = a$ .

```
"0 ∈ DPow(A) "
"a ∈ A ⇒ {a} ∈ DPow(A) "
```

The complement of a set  $X$  is defined by negating the formula used to define  $X$ . Intersection is done by conjoining the defining formulae, using the renaming techniques developed in the previous section. Union is then trivial by de Morgan's laws.

```
"X ∈ DPow(A) ⇒ (A - X) ∈ DPow(A) "
"[[X ∈ DPow(A); Y ∈ DPow(A)]] ⇒ X Int Y ∈ DPow(A) "
"[[X ∈ DPow(A); Y ∈ DPow(A)]] ⇒ X Un Y ∈ DPow(A) "
```

Thus  $DPow$  coincides with  $Pow$  (the real powerset operator) for finite sets.

```
"Finite(A) ⇒ DPow(A) = Pow(A) "
```



## 6.6. Proving that the ordinals are definable

In order to show that  $DPow$  is closed under other operations, we must be able to code their defining formulae as elements of the set *formula*. The treatment of the subset relation is typical. We begin by encoding the formula  $\forall z. z \in x \rightarrow z \in y$ . Below,  $x$  and  $y$  are de Bruijn indices, which are incremented to  $succ(x)$  and  $succ(y)$  because the quantifier introduces a new variable binding.

```
"subset_fm(x,y) ==
  Forall(Implies(Member(0,succ(x)), Member(0,succ(y))))"
```

The arguments are just de Bruijn indices because internalized formulae have no terms other than variables. It is trivial to prove that *subset\_fm* maps a pair of de Bruijn indices to a formula as follows.

```
"[x ∈ nat; y ∈ nat] ⇒ subset_fm(x,y) ∈ formula"
```

The arity of the formula is the maximum of those of its operands.

```
"[x ∈ nat; y ∈ nat] ⇒ arity(subset_fm(x,y)) = succ(x) ∪ succ(y)"
```

The following equivalence involves absoluteness, since it relates *subset\_fm* to the real subset relation,  $\subseteq$ . To reach this conclusion requires the additional assumption *Transset* ( $A$ ), saying that  $A$  is a transitive set. The premise  $x < length(env)$  puts a bound on  $x$  (which is a de Bruijn index), ensuring that  $nth(x, env)$  belongs to  $A$ .

```
"[x < length(env); y ∈ nat; env ∈ list(A); Transset(A)]
 ⇒ sats(A, subset_fm(x,y), env) ↔ nth(x,env) ⊆ nth(y,env)"
```

We must repeat this exercise (the details are omitted here) for the concepts of ‘transitive set’ and ‘ordinal’. This lets us prove that ordinals are definable, leading to a result involving ordinals and  $DPow$ .

```
"Transset(A) ⇒ {x ∈ A. Ord(x)} ∈ DPow(A)"
```

This lemma ultimately leads to a proof that  $\mathbf{L}$  contains all the ordinals.

## 6.7. Defining $\mathbf{L}$ , the constructible universe

The constant *Lset* formalizes the family of sets  $\{L_\alpha\}_{\alpha \in \mathbf{ON}}$ . Its definition in Isabelle/ZF uses a standard operator for transfinite recursion. We also define  $\mathbf{L} = \bigcup_{\alpha \in \mathbf{ON}} L_\alpha$ .

```
"Lset(i) == transrec(i, %x f. ⋃ y ∈ x. DPow(f'y))"
"L(x) == ∃ i. Ord(i) & x ∈ Lset(i)"
```

Some effort is required before we can transform the cryptic definition of *Lset* into the usual recursion equations. First, we prove Kunen’s [9, Lemma VI 1.6, p. 167], which states the transitivity and monotonicity of the  $L_\alpha$ .

```
"Transset(A) ⇒ Transset(DPow(A))"
"Transset(Lset(i))"
"i ≤ j → Lset(i) ⊆ Lset(j)"
```

Then we reach the 0, successor and limit equations for the  $L_\alpha$ .

```
"Lset(0) = 0"
"Lset(succ(i)) = DPow(Lset(i))"
"Limit(i) ⇒ Lset(i) = (⋃ y ∈ i. Lset(y))"
```

The basic properties of  $\mathbf{L}$ , as presented in Kunen's [9, IV, Section 1], are not hard to prove. For example,  $\mathbf{L}$  contains the ordinals.

"Ord( $i$ )  $\implies i \in Lset(succ(i))$  "

"Ord( $i$ )  $\implies L(i)$  "

## 6.8. Eliminating the arity function

The function *arity* can be surprisingly hard to reason about, particularly when we try to encode higher-order operators. Once we have established the basic properties of  $\mathbf{L}$ , we can prove its equivalence to a new definition that does not involve arities. Here is another form of definable powerset.

$$\begin{aligned} "DPow'(A) == \{ & X \in Pow(A). \\ & \exists env \in list(A). \exists p \in formula. \\ & X = \{x \in A. sats(A, p, Cons(x, env))\} \} " \end{aligned}$$

This version omits the constraint  $arity(p) \leq succ(length(env))$ , but is otherwise identical to  $DPow$ . The point is that if the environment is too short, attempted variable lookups will yield zero; recall the properties of  $nth$  from Section 6.2. If the set  $A$  is transitive, then it contains zero as an element. So the too-short environment can be padded to the right with zeroes.

"Transset( $A$ )  $\implies DPow(A) = DPow'(A)$  "

Each  $Lset(i)$  is a transitive set, so they can be expressed using  $DPow'$  rather than  $DPow$ , as follows.

"Lset( $i$ ) = transrec( $i$ ,  $\%x f. \bigcup_{y \in x. DPow'(f \text{ ` } y)$ )"

The equation above, proved by transfinite induction, lets us relativize  $Lset$  without having to formalize the functions *arity* and *length*. That eliminates a lot of work.

The following lemma is helpful for proving instances of separation. The first, quantified, premise asks for an equivalence between the real formula  $P$  and the internalized formula  $p$ . Often we can derive  $p$  from  $P$  automatically by supplying a set of suitable inference rules.

" $\llbracket \forall x \in Lset(i). P(x) \longleftrightarrow sats(Lset(i), p, Cons(x, env)) ;$   
 $env \in list(Lset(i)); p \in formula \rrbracket$   
 $\implies \{x \in Lset(i). P(x)\} \in DPow(Lset(i))$  "

Also, the lemma makes no reference to *arity*, thanks to the equivalence between  $DPow'$  and  $DPow$ .

## 6.9. The Zermelo–Fraenkel axioms hold in $\mathbf{L}$

Following [9, VI, Section 2], it is possible to prove that  $\mathbf{L}$  satisfies the Zermelo–Fraenkel axioms. Separation is the most difficult case, and is considered later.

### 6.9.1. Basic properties of $\mathbf{L}$

We begin with simple closure properties. Many of them involve exhibiting an element of *formula* describing the required set. We typically begin by starting in  $Lset(i)$  and proving that the required set belongs to  $Lset(succ(i))$ .

**L** is closed under unions.

$$"X \in \text{Lset}(i) \implies \text{Union}(X) \in \text{Lset}(\text{succ}(i))"$$

$$"L(X) \implies L(\text{Union}(X))"$$

**L** is closed under unordered pairs. More work is necessary because the sets  $a$  and  $b$  may be introduced at different ordinals.

$$"a \in \text{Lset}(i) \implies \{a\} \in \text{Lset}(\text{succ}(i))"$$

$$"\langle a, b \rangle \in \text{Lset}(i) \implies \{a, b\} \in \text{Lset}(\text{succ}(i))"$$

$$"\langle a, b \rangle \in \text{Lset}(i) \implies \{a, b\} \in \text{Lset}(i)"$$

$$"\langle L(a), L(b) \rangle \implies L(\{a, b\})"$$

Also,  $L_\alpha$  is closed under ordered pairs, provided that  $\alpha$  is a limit ordinal. This result is needed in order to apply the reflection theorem to **L**. Specifically, it is needed because my version of the reflection theorem [18] uses ordered pairs to cope with the possibility of a formula having any number of free variables.

$$"\langle a, b \rangle \in \text{Lset}(i) \implies \langle a, b \rangle \in \text{Lset}(\text{succ}(\text{succ}(i)))"$$

$$"\langle a, b \rangle \in \text{Lset}(i) \implies \langle a, b \rangle \in \text{Lset}(i)"$$

### 6.9.2. A rank function for **L**

Some proofs require the **L**-rank operator. Kunen [9, VI, 1.7] defines  $\rho(x)$  to denote the least  $\alpha$  such that  $x \in L_{\alpha+1}$ .

$$|\text{rank}(x) == \mu i. x \in \text{Lset}(\text{succ}(i))"$$

Here is one consequence of this definition.

$$"\text{Ord}(i) \implies x \in \text{Lset}(i) \iff L(x) \ \& \ \text{rank}(x) < i"$$

A more important result, whose proof involves  $\text{rank}$ , states that every set of constructible sets is included in some  $\text{Lset}$ .

$$"( \forall x \in A. L(x) ) \implies \exists i. \text{Ord}(i) \ \& \ A \subseteq \text{Lset}(i)"$$

This theorem is useful in proving that **L** satisfies the separation axiom. However, note that  $A \subseteq \mathbf{L}$  does not imply that  $A \in \mathbf{L}$ , not even if  $A$  is a set of natural numbers.

The  $\text{rank}$  operator is useful for proving that  $L$  satisfies the powerset axiom.

$$"L(X) \implies L(\{y \in \text{Pow}(X). L(y)\})"$$

Note that the powerset of  $X$  in  $L$  comprises all subsets of  $X$  that belong to  $L$ . It is potentially a superset of  $\text{DPow}(X)$ .

The  $\text{rank}$  operator also assists in the proof that  $L$  satisfies the replacement axiom. The idea is to use replacement on the ranks of the members of  $L$ .

$$"\langle L(X), \text{univalent}(L, X, Q) \rangle"$$

$$\implies \exists Y. L(Y) \ \& \ \text{Replace}(X, \lambda x y. Q(x, y) \ \& \ L(y)) \subseteq Y"$$

The proof of replacement is schematic, and therefore independent of the formula  $Q$ . However, it is the weak form of replacement. It concludes that the range of  $Q$  (viewed as a class function) is included in some constructible set  $Y$ . Strong replacement, which is the version that we really want, asserts that the range itself is constructible. Each instance of strong replacement requires proving an instance of the axiom of separation.

### 6.9.3. Instantiating the locale $M\_trivial$

Now we are ready to show that  $L$  satisfies all the properties that we assumed of the class  $M$ , which we used to develop the general theory of absoluteness. The class  $L$  is transitive.

$$\llbracket y \in x; L(x) \rrbracket \implies L(y)$$

The facts about  $L$  proved above can be summarized using the following relativized forms of the ZF axioms.

$$\begin{aligned} & \text{"Union\_ax}(L) \\ & \text{"upair\_ax}(L) \\ & \text{"power\_ax}(L) \\ & \text{"replacement}(L, P) \end{aligned}$$

We do not need  $L$  to satisfy the foundation axiom. However, this fact is a trivial consequence of the foundation axiom.

$$\text{"foundation\_ax}(L)$$

The theorems above are all that we need to prove  $\text{"PROP } M\_trivial(L) \text{"}$ . This theorem makes all the results proved in locale  $\text{"M\_trivial"}$  available as theorems about  $L$ . In particular, the absoluteness and closure results listed in Section 4.5 above apply to  $L$ .

## 7. Comprehension in $\mathbf{L}$

It remains to show that  $\mathbf{L}$  satisfies the axiom of separation. The proof requires the reflection theorem. As described elsewhere [18], my formalization of that theorem applies to any class  $\mathbf{M} = \bigcup_{\alpha \in \mathbf{ON}} M_\alpha$ , where the family  $\{M_\alpha\}_{\alpha \in \mathbf{ON}}$  is increasing and continuous. An additional condition is that if  $\alpha$  is a limit ordinal, then  $M_\alpha$  must be closed under ordered pairing. Isabelle's locale mechanism captures these requirements, and we can now instantiate the locale with the class  $\mathbf{L} = \bigcup_{\alpha \in \mathbf{ON}} L_\alpha$ . However, making it ready for practical use requires additional work.

### 7.1. The reflection relation

The reflection theorem states that if  $\phi(x_1, \dots, x_n)$  is a formula in  $n$  variables, then there exists a closed and unbounded class  $\mathbf{C}$  such that for all  $\alpha \in \mathbf{C}$  and  $x_1, \dots, x_n \in M_\alpha$  we have

$$\phi^{\mathbf{M}}(x_1, \dots, x_n) \iff \phi^{M_\alpha}(x_1, \dots, x_n).$$

In fact, we need only the weaker conclusion that  $\mathbf{C}$  is unbounded, which enables us to find a suitable  $\alpha > \beta$ , given any ordinal  $\beta$ .

Applying the reflection theorem yields an Isabelle formula describing the class  $\mathbf{C}$ . These formulae may be interesting in the case of small examples [18], but in typical applications they are huge. The trivial proofs, which merely refer to other instances of reflection, take minutes of computer time; the resulting theorems amount to pages of text. The obvious solution is to express the reflection theorem using an existential quantifier, but classes cannot be quantified over: they are formulae.

Fortunately, Isabelle makes a distinction between the object-logic (here first-order logic) and the metalogic (a fragment of higher-order logic) [12]. I was able to formalize a meta-existential quantifier. It lies outside of first-order logic – in particular, Isabelle will reject any attempt to use it in comprehensions. However, it can be used in top-level assertions, which is all that we need.

We can now define the reflection relation between two formulae  $P$  and  $Q$ .

```
"REFLECTS[P,Q] ==
  (??C. Closed_Unbounded(C) &
    (∀a. C(a) → (∀x ∈ Lset(a). P(x) ↔ Q(a,x)))) "
```

It relates the formulae just if there exists a class  $C$  satisfying the conclusion of the reflection theorem [18]. That is,  $C$  is a closed, unbounded class of ordinals  $\alpha$  such that  $P$  and  $Q$  agree on  $L_\alpha$ . The existential quantifier,  $??C$ , hides the prohibitively large formula describing this class. The following lemma illustrates the use of the reflection relation. Note that the quantification over classes has disappeared.

```
"[REFLECTS[P,Q]; Ord(i)]
  ⇒ ∃j. i < j & (∀x ∈ Lset(j). P(x) ↔ Q(j,x)) "
```

If  $REFLECTS[P, Q]$  holds, and  $i$  is an ordinal, then there exists a larger ordinal  $j$  for which  $P$  and  $Q$  agree. Our choice of  $i$  can make  $j$  arbitrarily large.

The general form of the reflection theorem uses the relativization operator, which cannot be expressed in Isabelle/ZF. However, given a specific formula  $\phi$ , we can generate an instance of the reflection theorem relating  $\phi^L$  and  $\phi^{L_\alpha}$ . Here is the base case, where normally  $P$  should have the form  $x \in y$  or  $x = y$ .

```
"REFLECTS[P, λa x. P(x)] "
```

Reflection relationships can be formed over the propositional connectives: here negation, conjunction and biconditionals.

```
"REFLECTS[P,Q] ⇒ REFLECTS[λx. ~P(x), λa x. ~Q(a,x)] "
```

```
"[REFLECTS[P,Q]; REFLECTS[P',Q']]
  ⇒ REFLECTS[λx. P(x) ∧ P'(x), λa x. Q(a,x) ∧ Q'(a,x)] "
```

```
"[REFLECTS[P,Q]; REFLECTS[P',Q']]
  ⇒ REFLECTS[λx. P(x) ↔ P'(x), λa x. Q(a,x) ↔ Q'(a,x)] "
```

Reflection relationships can be formed over the quantifiers.

```
"REFLECTS[ λx. P(fst(x),snd(x)), λa x. Q(a,fst(x),snd(x)) ]
  ⇒ REFLECTS[λx. ∃z[L]. P(x,z), λa x. ∃z∈Lset(a). Q(a,x,z)] "
```

In the conclusion, a quantification over  $L$  is related to one over  $L_\alpha$ , as suggested by the general form of the reflection theorem. The premise uses the projection operators for ordered pairs to introduce the new variable,  $z$ ; syntactically,  $\lambda x. P(fst(x), snd(x))$  is a unary formula.

## 7.2. Internalized formulae for some set-theoretic concepts

Every operator or concept that is used in an instance of the axiom of separation must be internalized. If the defining formula is complicated, then writing the corresponding element of *formula* requires a manual (and error-prone) translation into de Bruijn notation. The Isabelle/ZF development of constructibility theory contains about a hundred such encodings. A typical example resembles that shown in Section 6.6 above for *subset\_fm*. The first to be internalized are elementary concepts such as the empty set, unordered and ordered pairs, unions, intersections, domain and range.

The union predicate was defined in Section 4.1 as

$$\forall z. z \in Y \leftrightarrow z \in A \vee z \in B.$$

In the corresponding formula, the variables  $x$ ,  $y$  and  $z$  range over de Bruijn indices.

```
"union_fm(x,y,z) ==
  Forall (Iff (Member (0, succ (z)),
    Or (Member (0, succ (x)), Member (0, succ (y)))) )"
```

As for *subset\_fm* above, we can prove that *union\_fm* yields an element of the set *formula*. The theorem about satisfaction now takes the following form.

```
"[x ∈ nat; y ∈ nat; z ∈ nat; env ∈ list(A)]
 ⇒ sats(A, union_fm(x,y,z), env) ↔
  union(**A, nth(x,env), nth(y,env), nth(z,env))"
```

Here, *\*\*A* is Isabelle syntax for the class given by the set  $A$ ; that is,  $\{x \mid x \in A\}$ . The theorem above simply expresses the equivalence between the relational formula *union* and *union\_fm*, which is its translation into an element of the set *formula*. Such equivalences are usually trivial: they simply relate two syntaxes for formulae. They do not express the equivalence between *union\_fm* and  $\cup$ , which would be an instance of absoluteness.

After internalizing a predicate such as *union*, it makes sense to prove its instance of the reflection theorem too, since both results will be needed when proving instances of separation.

```
"REFLECTS[λx. union(L, f(x), g(x), h(x)),
  λi x. union(**Lset(i), f(x), g(x), h(x))]"
```

Most reflection proofs are trivial two-line scripts, as follows.

1. Unfold the concept's definition (here, *union*).
2. Repeatedly apply existing reflection theorems.

Each predicate is internalized similarly. Parts of the declarations and proofs can be copied from those of another predicate. However, getting the definition right requires careful attention to the original first-order definition.

### 7.3. Higher-order syntax

Higher-order syntax is ubiquitous in naive set theory. In the union  $\bigcup_{x \in A} B(x)$ , the higher-order variable  $B$  represents an indexed family of sets. In the function abstraction  $\lambda_{x \in A} b(x)$ , the higher-order variable  $b$  represents the function's body. Isabelle/ZF additionally uses higher-order syntax to express many forms of recursion, and so forth. Although this syntax is indispensable, it is also illegitimate: formal set theory has no non-trivial terms, let alone higher-order ones. We must formalize the conventions governing higher-order syntax into the language of set theory.

Converting a higher-order operator such as  $\lambda_{x \in A}. b(x)$  into relational form yields a higher-order predicate. Among its arguments is a predicate *is\_b* that expresses the function body,  $b$ , in relational form. If *is\_b* is purely relational, then so is the definiens of *is\_lambda*.

```
"is_lambda(M, A, is_b, z) ==
  ∀p[M]. p ∈ z ↔
  (∃u[M]. ∃v[M]. u ∈ A & pair(M,u,v,p) & is_b(u,v))"
```

This definition states that  $z$  is a  $\lambda$ -abstraction, provided that: (i) its elements are ordered pairs that satisfy  $is\_b$ , and (ii) the first component of each pair belongs to  $A$ .

The following predicate expresses the information that  $is\_f$  represents the relational version of  $f$  for arguments ranging over  $A$ .

```
"Relation1(M,A,is_f,f) ==
  ∀x[M]. ∀y[M]. x∈A → is_f(x,y) ↔ y = f(x) "
```

This abbreviation, and similarly *Relation2*, and so on, are useful for expressing absolute-ness results. If  $is\_b$  is the relational equivalent of  $b$ , and if the class  $M$  contains each  $b(m)$  for  $m \in A$ , then  $is\_lambda(M,A,is\_b,z)$  is the relational version of  $\lambda x \in A. b(x)$ . Thus  $\lambda$ -abstraction is absolute.

```
"[Relation1(M,A,is_b,b); M(A); ∀m[M]. m∈A → M(b(m)); M(z)]
⇒ is_lambda(M,A,is_b,z) ↔ (z = λx∈A. b(x)) "
```

Showing that  $M$  is closed under  $\lambda$ -abstraction requires a separate instance of strong replacement for each  $b$ .

```
"[strong_replacement(M, λx y. x∈A & y = (x, b(x)));
  M(A); ∀m[M]. m∈A → M(b(m)) ]
⇒ M(λx∈A. b(x)) "
```

Internalizing  $is\_lambda$  is not completely straightforward. The predicate argument,  $is\_b$ , becomes a variable ranging over the set *formula*.

```
"lambda_fm(p,A,z) ==
  Forall(Iff(Member(0,succ(z)),
    Exists(Exists(And(Member(1,A#+3),
      And(pair_fm(1,0,2), p)))))) "
```

Given a formula and two de Bruijn indices, *lambda\_fm* yields another formula.

```
"[p ∈ formula; x ∈ nat; y ∈ nat] ⇒ lambda_fm(p,x,y) ∈ formula"
```

However, there is no binding mechanism for expressing predicates that take arguments or refer to local variables. The formula  $p$  must refer to its first argument using the de Bruijn index 1, and to its second using the index 0 (both to be increased in the usual way if  $p$  contains quantifiers). If we are lucky, then we can arrange matters such that the actual arguments have the right indices, and otherwise we can force the indices to agree by introducing quantifiers and equalities: in the internalization of  $\forall x. \forall y. x = a \wedge y = b \rightarrow p$ , the variable with de Bruijn index 1 will refer to  $a$ , and similarly the index 0 will refer to  $b$ . If  $p$  contains free references to other variables, their de Bruijn indices must be increased by 3 because  $p$  is inserted into a context enclosed by three quantifiers.

The satisfaction theorem for  $is\_lambda$  formalizes the remarks above.

```
lemma sats_lambda_fm:
  assumes is_b_iff_sats ∈
    "!!a0 a1 a2.
      [a0∈A; a1∈A; a2∈A]
      ⇒ is_b(a1,a0) ↔ sats(A, p, Cons(a0,Cons(a1,Cons(a2,env)))) "
  shows
    "[x ∈ nat; y ∈ nat; env ∈ list(A)]
    ⇒ sats(A, lambda_fm(p,x,y), env) ↔
      is_lambda(**A, nth(x,env), is_b, nth(y,env)) "
```



The **assumes-shows** syntax eases the use of the complicated assumption which states that  $is\_b$  agrees with  $p$  for the fixed environment  $env$  extended with three additional elements of  $A$ . I have not been able to simplify the form of this theorem while retaining its generality.

It gets more complicated when one higher-order operator refers to another. One such operator has a quantifier nesting depth of 12. When an operator uses its higher-order argument more than once, we must ensure that the two contexts are similar, adding quantifiers if necessary to make the nesting depths agree.

Instances of the reflection theorem for higher-order operators must take into account the possibility of the higher-order argument's referring to local variables. Although  $is\_lambda$  expects  $is\_b$  to have only two arguments, below we formalize it with three arguments (plus its class argument). The extra argument is bound by the *REFLECTS* operator, allowing direct reference to elements of  $L$  or  $Lset(i)$ .

```
theorem is_lambda_reflection:
  assumes is_b_reflection:
    "!!f g h. REFLECTS[λx. is_b(L, f(x), g(x), h(x)),
      λi x. is_b(**Lset(i), f(x), g(x), h(x))]"
  shows "REFLECTS[λx. is_lambda(L, A(x), is_b(L,x), f(x)),
    λi x. is_lambda(**Lset(i), A(x), is_b(**Lset(i),x), f(x))]"
```

The arity of a higher-order function naturally depends upon that of its function argument. I found the properties so unintuitive, and their proofs so vexing, that I undertook the work described in Section 6.8, which eliminates the need for theorems concerning arities.

#### 7.4. Proving instances of separation

The set comprehension  $\{x \in A \mid \phi(x)\}$  comes from the separation axiom scheme instantiated to the formula  $\phi$ . The axiom of replacement yields a set that may be bigger than we want, again requiring an appeal to separation. Because I have not formalized the metatheory, the Isabelle/ZF development cannot express the proof that the separation scheme holds for  $\mathbf{L}$ . Each instance has to be proved individually. Fortunately, the proof scripts are nearly identical. Given  $\phi$ , the first step is to prove instance of the reflection theorem for that formula. The next step is to run a proof script corresponding to the sketch in [9, p. 169]. The formula  $\phi$  will of course be expressed using the relational language, using predicates such as *union*. Executing the proof script will automatically generate an internalized formula, with *union\_fm* in the corresponding place.

The lemmas outlined on the preceding pages suffice to prove many instances of separation. Consider the instance that justifies the existence of the intersection  $Inter(A)$ . We must first prove the corresponding instance of the reflection theorem.

```
"REFLECTS[λx. ∀y[L]. y∈A → x ∈ y,
  λi x. ∀y∈Lset(i). y∈A → x ∈ y]"
```

Such instances are written manually. A text editor can replace quantification over  $\mathbf{L}$  by quantification over  $L_\alpha$  in the second formula. The proof, almost always, is a one-line appeal to previous reflection theorems.

The statement of each instance of separation comes from the corresponding locale assumption. The locale refers to an arbitrary class  $\mathbf{M}$ , so we must replace  $\mathbf{M}$  by  $\mathbf{L}$ . The proof scripts are typically three lines long, and follow a regular pattern. Note that any parameters used in the separation formula (here  $A$ ) must be elements of  $\mathbf{L}$ .

```
"L(A) ⇒ separation(L, λx. ∀y[L]. y∈A → x∈y) "
```

The following instance of separation justifies relational composition. I leave the corresponding instance of reflection to the reader's imagination.

```
"[L(r); L(s)]
  => separation(L, λxz. ∃x[L]. ∃y[L]. ∃z[L]. ∃xy[L]. ∃yz[L].
    pair(L,x,z,xz) & pair(L,x,y,xy) & pair(L,y,z,yz) &
    xy∈s & yz∈r)"
```

After proving ten or so instances of separation, we arrive at the following cryptic theorem.

```
"PROP M_basic(L)"
```

This asserts that  $L$  satisfies the conditions of the locale  $M\_basic$ , namely all the instances of separation needed to derive well-founded recursion. The absoluteness and closure results proved in that locale (described in Section 4.6) now become applicable to  $L$ .

### 7.5. Automatic internalization of formulae

Isabelle's ability to translate formulae written in the relational language into members of *formula* simplifies the proofs of separation. Here is an example, from the proof of the instance shown above (about relational composition).

The first proof step applies a lemma for proving instances of separations. It yields a subgoal that has the assumptions  $r \in Lset(j)$  and  $s \in Lset(j)$ , where  $j$  is arbitrary. We have to prove that the comprehension belongs to the next level of the constructible hierarchy, namely  $DPow(Lset(j))$ .

```
{xz ∈ Lset(j) . ∃x∈Lset(j). ∃y∈Lset(j). ...} ∈DPow(Lset(j))
```

The second proof step applies a lemma for proving membership in  $DPow(Lset(j))$ . It yields three subgoals (see Figure 2). The first is to show the equivalence between the real formula

```
(∃xa∈Lset(j). ∃y∈Lset(j). ...)
```

and  $sats(Lset(j), ?p3(j), [x,r,s])$ . This is the satisfaction relation applied to  $?p3(j)$ , a 'logical variable' that can be replaced by any expression, possibly involving the bound variable  $j$ . The third subgoal in Figure 2, namely  $?p3(j) \in formula$ , checks that the chosen expression is an internalized formula. The second subgoal verifies that the environment,  $[r,s]$ , is well-typed – namely, that it belongs to  $list(Lset(j))$ .

The third proof step is as follows.

```
apply (rule sep_rules | simp)+
```

It applies some theorem of *sep\_rules*, then simplifies, then repeats if possible. This finishes the proof. All separation proofs have this form, save only that sometimes *sep\_rules* needs to be augmented with additional theorems.

Formula synthesis works in a way familiar to all Prolog programmers. Essentially, the theorems in *sep\_rules* comprise a Prolog program for generating internalized formulae. Most of the 'program clauses' relate real formulae to internal ones, and are derived from the basic properties of the satisfaction relation. For example, this one relates the real conjunction  $P \& Q$  with the term  $And(p,q)$ . The first two subgoals concern the synthesis of  $p$  and  $q$ . The third subgoal expresses a type constraint on *env*.

```
"[P ↔ sats(A,p,env); Q ↔ sats(A,q,env); env ∈ list(A)]
  => (P & Q) ↔ sats(A, And(p,q), env)"
```

1.  $\bigwedge j \ x. \llbracket L(r); L(s); r \in \text{Lset}(j); s \in \text{Lset}(j); x \in \text{Lset}(j) \rrbracket$   
 $\implies (\exists xa \in \text{Lset}(j). \exists y \in \text{Lset}(j). \exists z \in \text{Lset}(j). \text{pair}(**\text{Lset}(j), xa, z, x) \wedge (\exists xy \in \text{Lset}(j). \text{pair}(**\text{Lset}(j), xa, y, xy) \wedge (\exists yz \in \text{Lset}(j). \text{pair}(**\text{Lset}(j), y, z, yz) \wedge xy \in s \wedge yz \in r))) \longleftrightarrow \text{sats}(\text{Lset}(j), ?p3(j), [x, r, s]))$
2.  $\bigwedge j. \llbracket L(r); L(s); r \in \text{Lset}(j); s \in \text{Lset}(j) \rrbracket$   
 $\implies [r, s] \in \text{list}(\text{Lset}(j))$
3.  $\bigwedge j. \llbracket L(r); L(s); r \in \text{Lset}(j); s \in \text{Lset}(j) \rrbracket$   
 $\implies ?p3(j) \in \text{formula}$

Figure 2: Subgoals ready for automatic synthesis of a formula.

This ‘program clause’ relates the real quantification  $\forall x \in A. P(x)$  with the term  $\text{Forall}(p)$ . The first subgoal concerns the synthesis of  $p$  in an environment augmented with an arbitrary  $x \in A$ .

$\llbracket !!x. x \in A \implies P(x) \longleftrightarrow \text{sats}(A, p, \text{Cons}(x, \text{env})); \text{env} \in \text{list}(A) \rrbracket$   
 $\implies (\forall x \in A. P(x)) \longleftrightarrow \text{sats}(A, \text{Forall}(p), \text{env})$

The environment, which initially contains the parameters of the separation formula, gets longer with each nested quantifier. Each higher-order operator can add several elements to the environment, as mentioned above in Section 7.3.

A base case of synthesis relates the formula  $x \in y$  with the term  $\text{Member}(i, j)$ . The first two subgoals concern the synthesis of the de Bruijn indices  $i$  and  $j$ , as follows.

$\llbracket \text{nth}(i, \text{env}) = x; \text{nth}(j, \text{env}) = y; \text{env} \in \text{list}(A) \rrbracket$   
 $\implies (x \in y) \longleftrightarrow \text{sats}(A, \text{Member}(i, j), \text{env})$

Other base cases concern predicates of the relational language. This theorem, which relates the formula  $\text{union}(**A, x, y, z)$  with the term  $\text{union\_fm}(i, j, k)$ , is just a re-working of a theorem shown in Section 7.2 above.

$\llbracket \text{nth}(i, \text{env}) = x; \text{nth}(j, \text{env}) = y; \text{nth}(k, \text{env}) = z; i \in \text{nat}; j \in \text{nat}; k \in \text{nat}; \text{env} \in \text{list}(A) \rrbracket$   
 $\implies \text{union}(**A, x, y, z) \longleftrightarrow \text{sats}(A, \text{union\_fm}(i, j, k), \text{env})$

Given the subgoal  $\text{nth}(?i, \text{env}) = x$ , Isabelle can synthesize  $?i$ . This de Bruijn index is determined by  $x$ , which comes from the original formula, and  $\text{env}$ , which is given in advance. If  $x$  matches the head of the environment, then  $?i$  should be zero.

$\text{nth}(0, \text{Cons}(a, l)) = a$

If it does not match, however, then we should discard the head and attempt to synthesize a de Bruijn index using the tail.

$\llbracket \text{nth}(n, l) = x; n \in \text{nat} \rrbracket \implies \text{nth}(\text{succ}(n), \text{Cons}(a, l)) = x$

The automatic synthesis of internalized formulae saves much work in proofs of separation. In principle, we could rewrite every relational formula into its primitive constituents of membership and equality, removing the need for  $\text{union\_fm}$  and a hundred similar constants. However, if too few internalized primitives have been defined, formula synthesis takes many minutes.

## 8. Absoluteness of recursive datatypes

The Isabelle/ZF proofs discussed up to now include the construction of the class **L** and the proof that it is a model of the Zermelo–Fraenkel axioms. The next step is to show that **L** satisfies  $\mathbf{V} = \mathbf{L}$ . That fact follows by the absoluteness of constructibility, which follows by the absoluteness of satisfaction. Consulting the definition of *satisfies* reveals that we must still prove the absoluteness of lists, formulae, the function *nth*, and several other notions.

Isabelle/ZF defines the sets *list*(*A*) and *formula* automatically from their user-supplied descriptions [17]. These fixedpoint definitions have advantages, but their use of the powerset operator is an obstacle to proving absoluteness. For a start, *Pow*(*D*) must be eliminated from the following definition.

$$| \text{"lfp}(D, h) == \text{Inter}(\{X \in \text{Pow}(D) . h(X) \subseteq X\}) "$$

We proceed by formalizing standard concepts from domain theory [2, pp. 51–56]. A set is *directed* if it is non-empty and closed under least upper bounds. A function is *continuous* if it preserves the unions of directed sets.

$$| \begin{array}{l} \text{"directed}(A) == A \neq 0 \ \& \ (\forall x \in A. \ \forall y \in A. \ x \cup y \in A) " \\ \text{"contin}(h) == (\forall A. \ \text{directed}(A) \longrightarrow h(\bigcup A) = (\bigcup_{x \in A} h(x))) " \end{array}$$

We can prove that the least fixed point of a monotonic, continuous function *h* can be expressed as the union of the finite iterations of *h*.

$$| \llbracket \text{bnd\_mono}(D, h) ; \text{contin}(h) \rrbracket \Longrightarrow \text{lfp}(D, h) = (\bigcup_{n \in \text{nat}} h^n(0)) "$$

This equation eliminates not only *Pow*(*D*), but also every occurrence of *D*, which is the ‘bounding set’ [15, Section 2.2], and is itself typically defined using powersets.

In order to apply this equation, we must prove that standard datatype constructions preserve continuity. The case bases are that the constant function and the identity function are continuous.

$$\begin{array}{l} \text{"contin}(\lambda X. \ A) " \\ \text{"contin}(\lambda X. \ X) " \end{array}$$

Sums and products preserve continuity.

$$\begin{array}{l} \llbracket \text{contin}(F) ; \text{contin}(G) \rrbracket \Longrightarrow \text{contin}(\lambda X. \ F(X) + G(X)) " \\ \llbracket \text{contin}(F) ; \text{contin}(G) \rrbracket \Longrightarrow \text{contin}(\lambda X. \ F(X) \times G(X)) " \end{array}$$

These four lemmas cover all finitely branching datatypes, including lists and formulae.

### 8.1. Absoluteness for function iteration

In the equation given above for least fixed points, the term  $h^n(0)$  abbreviates *iterates*(*h*, *n*, 0). Isabelle/ZF defines *iterates*(*F*, *n*, *x*) by the obvious primitive recursion on  $n \in \text{nat}$ . The absoluteness of datatype definitions will follow from the absoluteness of *iterates*.

Recall that a well-founded function definition consists of a relation *r* and function body *H*; recall equation (1) of Section 5.2. Relativizing such a function definition requires relativizing *H* by an Isabelle/ZF relation, say *MH*. So to relativize *iterates*, we declare *is\_iterates* in terms of another predicate *iterates\_MH*, representing the body of the recursion.

```
"iterates_MH(M, isF, v, n, g, z) ==
  is_nat_case(M, v,
    λm u. ∃ gm[M]. fun_apply(M, g, m, gm) & isF(gm, u),
    n, z)"
"is_iterates(M, isF, v, n, Z) ==
  ∃ sn[M]. ∃ msn[M]. successor(M, n, sn) & membership(M, sn, msn) &
    is_wfrec(M, iterates_MH(M, isF, v), msn, n, Z)"
```

Incidentally,  $is\_nat\_case(M, a, isb, n, z)$  expresses case analysis on the natural number  $n$ . Note that we again work in the general setting of a class  $M$  satisfying certain conditions. Later, we shall prove that  $L$  meets those conditions.

The absoluteness theorem for well-founded recursion requires an instance of strong replacement for each function being defined. However,  $iterates$  is a higher-order function, and so technically  $iterates(F, n, x)$  involves a separate instance of well-founded recursion for each  $F$ . The function  $iterates\_replacement$  can express each required instance of replacement; its argument  $isF$  is the relational form of  $F$ .

```
"iterates_replacement(M, isF, v) ==
  ∀ n[M]. n ∈ nat →
    wfrec_replacement(M, iterates_MH(M, isF, v), Memrel(succ(n)))"
```

Assuming such an instance of replacement, and given that  $isF$  is the relational version of  $F$ , the absoluteness of  $iterates$  is a corollary of the general theorem about well-founded recursion.

```
"[iterates_replacement(M, isF, v); relation1(M, isF, F);
  n ∈ nat; M(v); M(z); ∀ x[M]. M(F(x))]]
⇒ is_iterates(M, isF, v, n, z) ↔ z = iterates(F, n, v)"
```

We similarly find that  $\mathbf{M}$  is closed under function iteration.

```
"[iterates_replacement(M, isF, v); relation1(M, isF, F);
  n ∈ nat; M(v); ∀ x[M]. M(F(x))]]
⇒ M(iterates(F, n, v))"
```

## 8.2. Absoluteness for lists and formulae

The formal treatment of continuity and  $iterates$  enables us to prove that lists and formulae are absolute.

The definition of lists generated by the Isabelle/ZF datatype [17] is too complicated to relativize easily. Instead, we prove its equivalence to a more abstract (and familiar) definition.

$$list(A) = lfp(univ(A), \lambda X. \{0\} + A * X)$$

The function given to  $lfp$  is continuous by construction, which lets us replace the least fixed point by iteration and eliminate the non-absolute set  $univ(A)$ .

$$contin(\lambda X. \{0\} + A * X) \\ list(A) = (\bigcup_{n \in nat.} (\lambda X. \{0\} + A * X)^n (0))$$

Now the absoluteness of  $list(A)$  is obvious. However, each element of this equation must be formalized in order to prove absoluteness. We begin by introducing an abbreviation for finite iterations of  $\lambda X. \{0\} + A * X$  – that is, for finite stages of the list construction.

$$list\_N(A, n) == (\lambda X. \{0\} + A * X)^n (0)$$

Next, we relativize the function  $\lambda X. \{0\} + A * X$ . The predicate *number1* recognizes the number 1, which equals the set  $\{0\}$ .

```
"is_list_functor(M,A,X,Z) ==
  ∃ n1[M]. ∃ AX[M].
    number1(M,n1) & cartprod(M,A,X,AX) & is_sum(M,n1,AX,Z) "
```

Next, we relativize the function *list\_N*, the finite iterations.

```
"is_list_N(M,A,n,Z) ==
  ∃ zero[M]. empty(M,zero) &
    is_iterates(M, is_list_functor(M,A), zero, n, Z) "
```

We relativize membership in *list*(*A*) as membership in *list\_N*(*A*, *n*) for some *n*. The predicate *finite\_ordinal* recognizes the natural numbers.

```
"mem_list(M,A,l) ==
  ∃ n[M]. ∃ listn[M].
    finite_ordinal(M,n) & is_list_N(M,A,n,listn) & l ∈ listn "
```

Finally, we can relativize the set of lists itself.

```
"is_list(M,A,Z) == ∀ l[M]. l ∈ Z ↔ mem_list(M,A,l) "
```

After proving the absoluteness of *list\_N*(*A*, *n*), we obtain the absoluteness of *list*(*A*), and we prove that *M* is closed under list formation.

```
"M(A) ⇒ M(list(A)) "
"⌈M(A); M(Z)⌋ ⇒ is_list(M,A,Z) ↔ Z = list(A) "
```

Formulae are proved absolute in just the same way. We express the set *formula* as an abstract least fixed point of a suitable function, prove that function to be continuous, and eliminate the *lfp* operator.

```
"formula = lfp(univ(0), λX. ((nat*nat) + (nat*nat)) + (X*X + X)) "
"contin(λX. ((nat*nat) + (nat*nat)) + (X*X + X)) "
"formula =
  (⋃ n ∈ nat. (λX. ((nat*nat) + (nat*nat)) + (X*X + X)) ^ n (0)) "
```

Proceeding as for lists, we define the predicates *is\_formula\_functor*, *is\_formula\_N*, *mem\_formula* and finally *is\_formula*. We obtain the desired theorems, as follows.

```
"M(formula) "
"M(Z) ⇒ is_formula(M,Z) ↔ Z = formula "
```

### 8.3. Recursion over lists and formulae

We have already seen (in Section 5) that functions defined by well-founded recursion are absolute. For mathematicians, that is enough to justify the absoluteness of functions defined recursively on lists or formulae. Proof tool users, however, must work through the details for each instance. Usually, automation makes it easy to apply general results to particular circumstances. However, the Isabelle/ZF translation of recursive function definitions is rather complicated (see [15, Sections 3.4 and 4.3.1].) There are good reasons for this complexity, such as support for a form of polymorphism. However, it makes the absoluteness proofs more difficult: the complications have to be taken apart and relativized, one by one.

At least there is no need to treat recursion over lists. Defining the class **L** involves only one list function, namely *nth*. Given a natural number *n* and a list *l*, this function returns

the  $n$ th element of  $l$ , counting from 0. Obviously, this amounts to taking the tail of the list  $n$  times and returning the head of the result. The recursion in  $nth$  is an instance of *iterates*.

Isabelle/ZF defines the head and tail functions  $hd$  and  $tl$ . The absoluteness proofs use modified versions called  $hd'$  and  $tl'$ , which extend  $hd$  and  $tl$  to return 0 if their argument is ill-formed (the details are unimportant). Relativization is simpler when a function's behaviour is fully specified. Now we can prove an equivalence for  $nth$ , as follows.

$$\llbracket xs \in list(A); n \in nat \rrbracket \implies nth(n, xs) = hd' (tl' ^ n (xs)) "$$

Its relational equivalent,  $is\_nth$ , has an obvious definition in terms of the relational equivalents of *iterates*,  $tl$  and  $hd$ .

$$\begin{array}{l} | \\ "is\_nth(M, n, l, Z) == \\ \quad \exists X[M]. is\_iterates(M, is\_tl(M), l, n, X) \ \& \ is\_hd(M, X, Z) " \end{array}$$

Absoluteness is proved with no effort.

$$\begin{array}{l} \llbracket M(A); n \in nat; l \in list(A); M(Z) \rrbracket \\ \implies is\_nth(M, n, l, Z) \longleftrightarrow Z = nth(n, l) " \end{array}$$

Recursion over lists is absolute in general. Proving this claim would require much work, and is unnecessary for proving that  $\mathbf{V} = \mathbf{L}$  is absolute. The function *satisfies* involves recursion over the datatype of formulae, and its absoluteness proof consists of several stages. Isabelle/ZF expresses recursion on datatypes in terms of  $\in$ -recursion, which is recursion on a set's rank [15, Section 3.4]. Absoluteness for  $\in$ -recursion will follow from that of well-founded recursion once we have established the absoluteness of  $\in$ -closure. Then we shall be in a position to consider recursion over formulae.

Five instances of strong replacement are necessary for the proofs sketched above. There are two each for the absoluteness of  $list(A)$  and *formula*, and one for the absoluteness of  $nth(n, l)$ . The locale  $M\_datatypes$  encapsulates these additional constraints on the class  $M$ . It is one of several locales used to keep track of instances of separation and replacement in this development.

#### 8.4. Absoluteness for $\in$ -closure

If  $A$  is a set, then its  $\in$ -closure is the smallest transitive set that includes  $A$ . Formally, the  $\in$ -closure of  $A$  is  $\bigcup_{n \in \omega} \bigcup^n(A)$ . Here,  $\bigcup^n(A)$  denotes the  $n$ -fold union of  $A$ , defined by  $\bigcup^0(A) = A$  and  $\bigcup^{m+1}(A) = \bigcup(\bigcup^m(A))$ . This is just another instance of *iterates*, which we can prove as follows.

$$eclose(A) = (\bigcup_{n \in nat}. Union^n(A)) "$$

Relativization proceeds as it did for lists. The details are omitted here, but they culminate in the definition of a relational version of  $eclose(A)$ .

$$\begin{array}{l} | \\ "is\_eclose(M, A, Z) == \forall u[M]. u \in Z \longleftrightarrow mem\_eclose(M, A, u) " \end{array}$$

The standard membership and absoluteness results follow.

$$\begin{array}{l} M(A) \implies M(eclose(A)) " \\ \llbracket M(A); M(Z) \rrbracket \implies is\_eclose(M, A, Z) \longleftrightarrow Z = eclose(A) " \end{array}$$



### 8.5. Absoluteness for *transrec*

The Isabelle/ZF operator *transrec* expresses  $\in$ -recursion, which includes transfinite recursion as a special case.

$$\text{transrec}(a, H) = H(a, \lambda_{x \in a}. \text{transrec}(x, H)).$$

Its definition is a straightforward combination of the operators *eclose*, *wfrec* (which expresses well-founded recursion), and *Memrel* (which encodes the membership relation as a set). Thus the definition of the relational version, *is\_transrec*, is also straightforward. Our previous results lead directly to a proof of absoluteness.

```
"[[transrec_replacement(M,MH,i); relativize2(M,MH,H);
  Ord(i); M(i); M(z);
  ∀x[M]. ∀g[M]. function(g) → M(H(x,g))]]
⇒ is_transrec(M,MH,i,z) ↔ z = transrec(i,H)"]
```

We similarly find that *M* is closed under  $\in$ -recursion.

```
"[[transrec_replacement(M,MH,i); relativize2(M,MH,H);
  Ord(i); M(i);
  ∀x[M]. ∀g[M]. function(g) → M(H(x,g))]]
⇒ M(transrec(i,H))"]
```

In these theorems, *transrec\_replacement* abbreviates a specific use of *wfrec\_replacement*, which justifies this particular recursive definition (recall Section 5.3).

### 8.6. Recursion over formulae

The Isabelle/ZF treatment of recursive functions on datatypes involves non-absolute concepts, namely the cumulative hierarchy  $\{V_\alpha\}_{\alpha \in \mathbf{ON}}$  and the rank function [15, Section 3.6]. For proving absoluteness, I proved an equation stating that recursion over formulae could be expressed differently. The new formulation refers to the *depth* of a formula, defined as follows.

```
"depth(Member(x,y)) = 0"
"depth(Equal(x,y)) = 0"
"depth(Nand(p,q)) = succ(depth(p) ∪ depth(q))"
"depth(Forall(p)) = succ(depth(p))"
```

Introducing *depth* seems to be a step backwards, since it requires relativizing another recursive function on formulae. However, we can express the depth of a formula in terms of *is\_formula\_N*, which we need anyway (Section 8.2); *is\_formula\_N*(*M*, *n*, *F*) holds just if *F* is the set of formulae generated by *n* unfoldings of the datatype definition – which is all formulae of depth less than *n*. A formula *p* has depth *n* if it satisfies *is\_formula\_N*(*M*, *succ*(*n*), *F*) and not *is\_formula\_N*(*M*, *n*, *F*).

```
"is_depth(M,p,n) ==
  ∃sn[M]. ∃formula_n[M]. ∃formula_sn[M].
    is_formula_N(M,n,formula_n) & p ∉ formula_n &
    successor(M,n,sn) &
    is_formula_N(M,sn,formula_sn) & p ∈ formula_sn"
```

Working from this definition, we find that the depth of a formula is absolute.

$$[[p \in \text{formula}; n \in \text{nat}]] \implies \text{is\_depth}(M,p,n) \longleftrightarrow n = \text{depth}(p)''$$

For relativization, I modified the standard Isabelle/ZF treatment of recursion over formulae, replacing the set  $V_\alpha$  by *formula* and the rank of a set by the depth of a formula. If  $f$  is a recursive function on formulae, then the evaluation of  $f(p)$  begins by determining the depth of  $p$ , say  $n$ . Then the recursion equation for  $f$  is unfolded  $n + 1$  times, using transfinite recursion. The resulting nonrecursive function is finally applied to  $p$ . This approach unfortunately needs an explicit  $\lambda$ -abstraction over formulae and another instance of the replacement axiom. With the benefit of hindsight, I might have saved much work by seeking simpler ways of expressing recursion over formulae, such as by well-founded recursion on the subformula relation.

The recursive definition of a function  $f$  is specified by four parameters  $a, b, c$  and  $d$ , corresponding to the four desired recursion equations.

$$\begin{aligned} f(\text{Member}(x, y)) &= a(x, y) \\ f(\text{Equal}(x, y)) &= b(x, y) \\ f(\text{Nand}(p, q)) &= c(p, f(p), q, f(q)) \\ f(\text{Forall}(p)) &= d(p, f(p)) \end{aligned}$$

Given the datatype definition of *formula*, Isabelle/ZF automatically defines the operator *formula\_rec* for expressing recursive functions. The term *formula\_rec*( $a, b, c, d, p$ ) denotes the value of the function  $f$  above applied to the argument  $p$ . More concisely, *formula\_rec*( $a, b, c, d$ ) denotes the function  $f$  itself. The details of the definitions are illustrated elsewhere, using the example of lists [15, Section 4.3].

In order to express the recursion theorem, it helps to have first defined an abbreviation for its case analysis on formulae.

$$\begin{aligned} &\text{"formula\_rec\_case}(a, b, c, d, h) == \\ &\quad \text{formula\_case } (a, b, \\ &\quad \quad \lambda u \ v. \ c(u, \ v, \ h \ ' \ \text{succ}(\text{depth}(u)) \ ' \ u, \\ &\quad \quad \quad h \ ' \ \text{succ}(\text{depth}(v)) \ ' \ v), \\ &\quad \lambda u. \ d(u, \ h \ ' \ \text{succ}(\text{depth}(u)) \ ' \ u)) \text{"} \end{aligned}$$

Now we can express recursion on formulae in terms of absolute concepts.

$$\begin{aligned} &\text{"}p \in \text{formula} \implies \\ &\quad \text{formula\_rec}(a, b, c, d, p) = \\ &\quad \text{transrec } (\text{succ}(\text{depth}(p)), \\ &\quad \quad \lambda x \ h. \ \text{Lambda}(\text{formula}, \text{formula\_rec\_case}(a, b, c, d, h))) \ ' \ p \text{"} \end{aligned}$$

The proof is by structural induction on  $p$ . Note that the argument  $h$  of *formula\_rec\_case* is a partially unfolded recursive function of two arguments. It is a *curried* function: given its first argument, a natural number, the result is another function that can accept a second argument, a formula. In *formula\_rec\_case*, the second argument is some formula  $u$  (which in practice will be a subformula of  $p$ ), and the first is  $\text{succ}(\text{depth}(u))$ . The intuition behind this theorem may be obscure, but that is no obstacle to proving absoluteness. Many routine details must be taken care of, including relativization and absoluteness for the formula constructors *Member*, *Equal*, *Nand* and *Forall*, and for the operator *formula\_case*.

Obviously, *formula\_rec* is a higher-order function. Its absoluteness proof depends upon absoluteness assumptions for the function arguments  $a, b, c$  and  $d$ . Its relational version needs those arguments to be expressed in relational form as predicates *is\_a*, *is\_b*, *is\_c* and *is\_d*. The absoluteness theorem depends upon ten assumptions in all: two each for *is\_a*, *is\_b*, *is\_c* and *is\_d*, and two instances of replacement.

After many intricate but uninteresting details, we arrive at two key theorems. If the class  $M$  is closed under the parameters  $a$ ,  $b$ ,  $c$  and  $d$ , then it is closed under the corresponding recursion.

$$"p \in \text{formula} \implies M(\text{formula\_rec}(a, b, c, d, p))"$$

Recursion over formulae is absolute.

$$\begin{aligned} & " \llbracket p \in \text{formula}; M(z) \rrbracket \\ & \implies \text{is\_formula\_rec}(M, MH, p, z) \longleftrightarrow z = \text{formula\_rec}(a, b, c, d, p) " \end{aligned}$$

In this theorem,  $MH$  abbreviates the relativization of the argument of  $\text{transrec}$  shown above.

$$\begin{aligned} & "MH(u::i, f, z) == \\ & \quad \forall fml[M]. \text{is\_formula}(M, fml) \longrightarrow \\ & \quad \quad \text{is\_lambda} \\ & \quad (M, fml, \text{is\_formula\_case}(M, \text{is\_a}, \text{is\_b}, \text{is\_c}(f), \text{is\_d}(f)), z) " \end{aligned}$$

## 9. Absoluteness for $\mathbf{L}$

In order to prove that  $\mathbf{V} = \mathbf{L}$ , we must prove the absoluteness of three main functions.

1. *satisfies*, the satisfaction function on formulae;
2.  $DPow$ , the definable powerset function;
3.  $Lset$ , which expresses the levels  $L_\alpha$  of the constructible hierarchy.

Of these functions,  $Lset$  is defined by transfinite recursion from  $DPow$ , which in turn has a straightforward definition in terms of *satisfies*. Proving the absoluteness of *satisfies* is, however, very complicated.

Absoluteness of *satisfies* is merely an instance of the absoluteness of recursion over formulae, and is therefore trivial. That does not relieve us of the task of formalizing the details. The file containing the *satisfies* absoluteness proof is one of the largest in the entire development. This file divides into two roughly equal parts.

The first half contains internalizations and reflection theorems for operators such as *depth* and *formula\_case*. It expresses the four cases of *satisfies* in both functional and relational form, and proves absoluteness for each case. Six instances of strong replacement are required: one for each case of the recursion (because each contains a  $\lambda$ -abstraction), another to justify the use of *transrec*, and yet another to justify the  $\lambda$ -abstraction in *formula\_rec*. These axioms are assumed to hold for an arbitrary class model  $M$ . They are used to show that the formalization satisfies the conditions of the absoluteness theorem for *formula\_rec* described in the previous section.

The second half of the file is devoted to proving that the six instances of replacement hold in  $L$ . The four cases of the recursion (in their relational form) must each be internalized. This tiresome task involves, as always, translating a definition involving real formulae into one using internalized formulae. Then the six instances of replacement are justified. Finally, the pieces are put together.

### 9.1. Proving that *satisfies* is absolute

Working in the class  $M$ , we assume additional instances of the replacement axiom and apply them to the definition of *satisfies*, which is reproduced here.

```
"satisfies(A, Member(x,y)) =
  (λenv ∈ list(A). bool_of_o (nth(x,env) ∈ nth(y,env))) "
"satisfies(A, Equal(x,y)) =
  (λenv ∈ list(A). bool_of_o (nth(x,env) = nth(y,env))) "
"satisfies(A, Nand(p,q)) =
  (λenv ∈ list(A). not ((satisfies(A,p) 'env) and
    (satisfies(A,q) 'env))) "
"satisfies(A, Forall(p)) =
  (λenv ∈ list(A). bool_of_o
    (∀x∈A. satisfies(A,p) ' (Cons(x,env)) = 1)) "
```

Many additional concepts must be internalized. Consider the predicate *is\_depth*, which formalizes the depth of a formula.

```
"depth_fm(p,n) ==
  Exists(Exists(Exists(
    And(formula_N_fm(n#+3,1),
    And(Neg(Member(p#+3,1)),
    And(succ_fm(n#+3,2),
    And(formula_N_fm(2,0), Member(p#+3,0))))))) "
```

We prove the usual theorem relating the satisfaction of *depth\_fm* to the truth of *is\_depth*.

```
"[x ∈ nat; y < length(env); env ∈ list(A)]
⇒ sats(A, depth_fm(x,y), env) ↔
  is_depth(**A, nth(x,env), nth(y,env)) "
```

Then we generate yet another instance of the reflection theorem.

```
"REFLECTS[λx. is_depth(L, f(x), g(x)),
  λi x. is_depth(**Lset(i), f(x), g(x))]"
```

The internalization of *is\_formula\_case* is omitted here, but its definition is fifteen lines long and contains eleven quantifiers. The theorem statements relating *is\_formula\_case* to *formula\_case* are also long and complicated. They are of course also higher-order, requiring the methods of Section 7.3.

In order to relativize *satisfies*, we must first define constants corresponding to *formula\_rec*'s parameters *a*, *b*, *c* and *d*. Here are the two base cases.

```
"satisfies_a(A) ==
  λx y. λenv ∈ list(A). bool_of_o (nth(x,env) ∈ nth(y,env)) "
"satisfies_b(A) ==
  λx y. λenv ∈ list(A). bool_of_o (nth(x,env) = nth(y,env)) "
```

In the two recursive cases, the variables *rp* and *rq* denote the values returned on the recursive calls for *p* and *q*, respectively.

```
"satisfies_c(A) ==
  λp q rp rq. λenv ∈ list(A). not(rp ' env and rq ' env) "
"satisfies_d(A) ==
  λp rp. λenv ∈ list(A). bool_of_o (∀x∈A. rp ' (Cons(x,env)) = 1) "
```

Each of these functions is then re-expressed in relational form. Here is the first.

```
"satisfies_is_a(M,A) ==
  λx y zz. ∀ lA[M]. is_list(M,A,lA) →
    is_lambda(M, lA,
      λenv z. is_bool_of_o(M,
        ∃ nx[M]. ∃ ny[M].
          is_nth(M,x,env,nx) & is_nth(M,y,env,ny) & nx ∈ ny, z),
      zz) "
```

Once we have done the other three, we can define an instance of *MH* for *satisfies*, expressing the body of the recursion as a predicate.

```
"satisfies_MH ==
  λM A u f z.
    ∀ fml[M]. is_formula(M, fml) →
      is_lambda (M, fml,
        is_formula_case (M, satisfies_is_a(M, A),
                          satisfies_is_b(M, A),
                          satisfies_is_c(M, A, f),
                          satisfies_is_d(M, A, f)),
        z) "
```

Finally, *satisfies* can itself be relativized.

```
"is_satisfies(M, A) == is_formula_rec (M, satisfies_MH(M, A)) "
```

This lemma relates the fragments defined above to the original primitive recursion in *satisfies*. Induction is not required: the definitions are directly equal!

```
"satisfies(A, p) =
  formula_rec (satisfies_a(A), satisfies_b(A),
    satisfies_c(A), satisfies_d(A), p) "
```

At this point we must assume (by declaring a locale) the six instances of replacement mentioned above. That enables us to prove absoluteness for the parameters *a*, *b*, *c* and *d* used to define *satisfies*. For example, the class *M* is closed under *satisfies\_a*.

```
"[[M(A); x ∈ nat; y ∈ nat]] ⇒ M(satisfies_a(A, x, y)) "
```

This theorem states that *satisfies\_is\_a*(*M*, *A*, *x*, *y*, *zz*) is the relational equivalent of *satisfies\_a*(*A*, *x*, *y*), provided that *x* and *y* belong to the set *nat*.

```
"M(A) ⇒
  Relation2(M, nat, nat, satisfies_is_a(M, A), satisfies_a(A)) "
```

It can be seen as an absoluteness result, subject to typing conditions on *x* and *y*. Proofs are obviously easier if the absoluteness results are unconditional, but sometimes typing conditions are difficult to avoid.

Analogous theorems are proved for *satisfies\_is\_b*, *satisfies\_is\_c* and *satisfies\_is\_d*. Thus we use the first four instances of replacement. The last two instances, which are specific to *satisfies*, allow us to discharge the more general instances of replacement that are conditions of *formula\_rec*'s absoluteness theorem. We ultimately obtain absoluteness for *satisfies* as follows.

```
"[[M(A); M(z); p ∈ formula]]
⇒ is_satisfies(M, A, p, z) ↔ z = satisfies(A, p) "
```

## 9.2. Proving the instances of replacement for **L**

Now we must justify those six instances of strong replacement by proving that they hold in **L**. Recall that strong replacement is the conjunction of replacement (which holds schematically in **L**, but may yield too big a set) and an appropriate instance of separation (Section 4.3).

As always, proving instances of separation requires the internalization of many formulae. Isabelle can do this automatically, but unless it is given enough internalized formulae to use

as building blocks, the translation requires much time and space. I internalized many concepts manually, declaring their internal counterparts as constants, and proving their correspondence with the original concepts. Here is the internal equivalent of *satisfies\_is\_a*.

```
"satisfies_is_a_fm(A,x,y,z) ==
  Forall(
    Implies(is_list_fm(succ(A),0),
      lambda_fm(
        bool_of_o_fm(Exists(
          Exists(And(nth_fm(x#+6,3,1),
                    And(nth_fm(y#+6,3,0),
                      Member(1,0))))),0),
        0, succ(z))))"
```

Obviously, the same task must be undertaken for the other *satisfies* relations, and for the concepts used in their definitions. Finally, we can internalize the body of *satisfies* as follows.

```
"satisfies_MH_fm(A,u,f,zz) ==
  Forall(
    Implies(is_formula_fm(0),
      lambda_fm(
        formula_case_fm(satisfies_is_a_fm(A#+7,2,1,0),
                        satisfies_is_b_fm(A#+7,2,1,0),
                        satisfies_is_c_fm(A#+7,f#+7,2,1,0),
                        satisfies_is_d_fm(A#+6,f#+6,1,0),
                        1, 0),
        0, succ(zz))))"
```

Now, we can prove the six instances of replacement. Here is the first one, for the *Member* case of *satisfies*.

```
"[L(A); x ∈ nat; y ∈ nat]
⇒ strong_replacement
  (L, λenv z. ∃bo[L]. ∃nx[L]. ∃ny[L].
    env ∈ list(A) & is_nth(L,x,env,nx) & is_nth(L,y,env,ny) &
    is_bool_of_o(L, nx ∈ ny, bo) &
    pair(L, env, bo, z))"
```

The theorem statement may look big, but the proof has only four commands. The corresponding instance of the reflection theorem (not shown) is twice as big, but its proof has only one command.

We proceed to prove the fifth instance of replacement.

```
"[n ∈ nat; L(A)] ⇒ transrec_replacement(L, satisfies_MH(L,A), n)"
```

Finally, we prove the sixth instance of replacement.

```
"[L(g); L(A)] ⇒
strong_replacement (L,
  λx y. mem_formula(L,x) &
    (∃c[L]. is_formula_case(L, satisfies_is_a(L,A),
                        satisfies_is_b(L,A),
                        satisfies_is_c(L,A,g),
                        satisfies_is_d(L,A,g), x, c) &
    pair(L, x, c, y)))"
```

Our reward for this huge effort is that the absoluteness of *satisfies* now holds for *L*.

```
"[L(A); L(z); p ∈ formula]
⇒ is_satisfies(L,A,p,z) ↔ z = satisfies(A,p)"
```

### 9.3. Absoluteness of the definable powerset

Conceptually, the absoluteness of  $DPow$  is trivial, since it is just a comprehension involving *satisfies*. The formal details require a modest effort. There are more internalizations, such as that of *is\_formula\_rec*. Note that concepts have to be internalized only if they appear in an instance of separation, which may happen only long after the concept is first relativized. Unfortunately, *formula\_rec* is a complex higher-order function; in its relational form, one argument gets enclosed within eleven quantifiers. Completing this task enables us to internalize *is\_satisfies*.

```
"satisfies_fm(x) == formula_rec_fm(satisfies_MH_fm(x#+5#+6,2,1,0))"
```

Recall that  $DPow$  is the definable powerset operator. It has a variant form,  $DPow'$ , that does not involve the function *arity*. The two operators agree on transitive sets, so in particular we can use  $DPow'$  to construct  $L$ . Now we must relativize  $DPow'$ . Its definition refers to the powerset operator, which is not absolute. It can equivalently be expressed using a set comprehension, which here represents an appeal to the replacement axiom.

```
"DPow'(A) = {z . ep ∈ list(A) × formula,
               ∃ env ∈ list(A). ∃ p ∈ formula.
               ep = (env,p) &
               z = {x∈A. sats(A, p, Cons(x,env))}}"
```

Within the comprehension is another comprehension, which appeals to separation. The formula *sats*( $A, p, Cons(x,env)$ ) needs to be relativized (as the predicate *is\_DPow\_sats*) and internalized. Then we again extend the list of assumptions about the class  $M$  to include these instances of replacement and separation. Using them, we can prove that  $M$  is closed under definable powersets.

```
"M(A) ⇒ M(DPow'(A))"
```

We can also express the equation for  $DPow'$  shown above in relational form, defining the predicate *is\_DPow'*, and prove absoluteness.

```
"[M(A); M(Z)] ⇒ is_DPow'(M,A,Z) ↔ Z = DPow'(A)"
```

To make these results available for  $L$ , we must first prove that  $L$  satisfies the new instances of replacement and separation. Here is the latter.

```
"[L(A); env ∈ list(A); p ∈ formula]
 ⇒ separation(L, λx. is_DPow_sats(L,A,env,p,x))"
```

### 9.4. Absoluteness of constructibility

The proof that  $L$  satisfies  $V = L$  is nearly finished. Only the operator  $Lset$ , which denotes the levels of the constructible hierarchy, remains to be proved absolute. Recall that it can be expressed using  $DPow'$  as follows.

```
"Lset(i) = transrec(i, %x f. ⋃ y∈x. DPow'(f ` y))"
```

So now we must internalize the predicate *is\_DPow'*. First we must internalize the operators used in its definition. Among those are the predicate *is\_Collect*, which recognizes set comprehensions. The equation for  $Lset$  above involves two further instances of replacement: one for the use of *transrec*, and another for the indexed union. Adding them to our list of constraints on  $M$  allows us to prove that that class is closed under the  $Lset$  operator.

" $\llbracket \text{Ord}(i); M(i) \rrbracket \implies M(\text{Lset}(i))$ "

We can also define its relational version as follows.

" $\text{is\_Lset}(M, a, z) ==$   
 $\text{is\_transrec}(M, \lambda x f u. u = (\bigcup_{y \in x}. \text{DPow}'(f \restriction y)), a, z)$ "

Notice that this definition is not purely relational. That is all right because  $\text{is\_Lset}$  is not used in any instance of separation, and thus need not be internalized. We can now prove that the constructible hierarchy is absolute.

" $\llbracket \text{Ord}(i); M(i); M(z) \rrbracket \implies \text{is\_Lset}(M, i, z) \longleftrightarrow z = \text{Lset}(i)$ "

As remarked earlier, results such as this express absoluteness because the class model  $M$  drops out of the right-hand side. The left-hand side refers to our formalization of  $L_\alpha$  in  $M$ , which by the theorem is equivalent to  $L_\alpha$  itself. As always, making this result available to  $L$  requires proving the new instances of replacement. I omit the details, which contain nothing instructive.

We can finally formalize  $\mathbf{L}^M$ , the relativization of  $\mathbf{L}$ . A set  $x$  is constructible (with respect to any class  $M$  satisfying the specified ZF axioms) provided that there exists an ordinal  $i$  and a level of the constructible hierarchy  $Li$  such that  $x \in Li$ .

" $\text{constructible}(M, x) ==$   
 $\exists i[M]. \exists Li[M]. \text{ordinal}(M, i) \ \& \ \text{is\_Lset}(M, i, Li) \ \& \ x \in Li$ "

The following theorem is a trivial consequence of the absoluteness results and the definitions of  $\text{constructible}$  and  $L$ .

" $L(x) \implies \text{constructible}(L, x)$ "

This theorem expresses our goal, namely that  $\mathbf{V} = \mathbf{L}$  holds in  $\mathbf{L}$  or, more formally,  $(\mathbf{V} = \mathbf{L})^{\mathbf{L}}$ , for this statement is equivalent to  $(\forall x. \mathbf{L}(x))^{\mathbf{L}}$  and thus to  $\forall x. \mathbf{L}(x) \rightarrow \mathbf{L}^{\mathbf{L}}(x)$ . We can drop the universal quantifier. The antecedent of the implication is formalized as  $L(x)$ , and the consequent as  $\text{constructible}(L, x)$ . This proof ends the most difficult part of the development.

## 10. The axiom of choice in $\mathbf{L}$

The formalization confirms that  $\mathbf{V} = \mathbf{L}$  is consistent with the axioms of set theory. Obviously any consequence of  $\mathbf{V} = \mathbf{L}$ , such as the axiom of choice, is consistent with those axioms too. Proving consequences of  $\mathbf{V} = \mathbf{L}$  involves working in an entirely different way, and a much pleasanter one. Dispensing with the relational language, relativization, internalization and absoluteness, we can instead work in native set theory with the additional axiom  $\mathbf{V} = \mathbf{L}$ .

Assuming that  $\mathbf{V} = \mathbf{L}$ , the proof of the axiom of choice is simple [9, p. 173]. It suffices to prove that every set can be well-ordered. In fact, we can well-order the whole of  $\mathbf{L}$ . The set of internalized formulae is countable, and therefore well-ordered. The well-ordering of  $\mathbf{L}$  derives from its cumulative construction, and from the well-ordering of formulae. For  $x, y \in \mathbf{L}$ , say that  $x$  precedes  $y$  if

- $x$  originates earlier than  $y$  in the constructible hierarchy – that is, there is some  $\alpha$  such that  $x \in L_\alpha$  and  $y \notin L_\alpha$ ; or
- $x$  and  $y$  originate at the same level  $L_\alpha$ , but the combination of defining formula and parameters for  $x$  lexicographically precedes the corresponding combination for  $y$ .



Each element of  $L_{\alpha+1}$  is a subset of  $L_\alpha$  that can be defined by a formula, possibly involving parameters from  $L_\alpha$ . We can assume the induction hypothesis that  $L_\alpha$  is well-ordered. Before we can undertake this transfinite induction, we must complete several tasks. We must

1. exhibit a well-ordering on lists, for the parameters of a definable subset;
2. exhibit a well-ordering on formulae;
3. combine these to obtain a well-ordering of the definable powerset;
4. show how to extend our well-ordering to the limit case of the transfinite induction.

### 10.1. A well-ordering for lists

First we inductively define a relation on lists: the lexicographic extension of a relation on the list's elements. Let  $r$  denote a relation over the set  $A$ . Then the relation  $rlist(A, r)$  is the least set closed under the following rules.

$$\begin{aligned} & \llbracket \text{length}(l') < \text{length}(l); l' \in list(A); l \in list(A) \rrbracket \\ & \implies \langle l', l \rangle \in rlist(A, r) \\ \\ & \llbracket \langle l', l \rangle \in rlist(A, r); a \in A \rrbracket \\ & \implies \langle Cons(a, l'), Cons(a, l) \rangle \in rlist(A, r) \\ \\ & \llbracket \text{length}(l') = \text{length}(l); \langle a', a \rangle \in r; \\ & \quad l' \in list(A); l \in list(A); a' \in A; a \in A \rrbracket \\ & \implies \langle Cons(a', l'), Cons(a, l) \rangle \in rlist(A, r) \end{aligned}$$

Informally, the list  $l'$  precedes another list  $l$  if

1.  $l'$  is shorter than  $l$ , or
2. the lists have the same head and the tail of  $l'$  precedes that of  $l$ , or
3. the lists have the same length and the head of  $l'$  precedes that of  $l$  under the ordering on list elements.

If the element ordering is linear, then so is the list ordering. This theorem has a fourteen-line proof script involving a double structural induction on lists.

$$linear(A, r) \implies linear(list(A), rlist(A, r))$$

If the element ordering is well-founded, then so is the list ordering. This theorem is proved by induction on the length of the list followed by inductions over the element ordering and the list ordering. The proof script is under twenty lines, but the argument is complicated.

$$well\_ord(A, r) \implies well\_ord(list(A), rlist(A, r))$$

### 10.2. A well-ordering on formulae

Gödel-numbering is the obvious way to well-order the set of formulae. An injection from the set of formulae into the set of natural numbers is easily defined by recursion on the structure of formulae. However, it requires an injection from pairs of natural numbers to natural numbers. The enumeration function for formulae takes this injection as its first argument,  $f$ , as follows.

$$\begin{aligned} enum(f, Member(x, y)) &= f \cdot \langle 0, f \cdot \langle x, y \rangle \rangle \\ enum(f, Equal(x, y)) &= f \cdot \langle 1, f \cdot \langle x, y \rangle \rangle \\ enum(f, Nand(p, q)) &= f \cdot \langle 2, f \cdot \langle enum(f, p), enum(f, q) \rangle \rangle \\ enum(f, Forall(p)) &= f \cdot \langle succ(2), enum(f, p) \rangle \end{aligned}$$

There are several well-known injections from  $\omega \times \omega$  into  $\omega$ , but defining one of them and proving it to be injective would involve some effort. Instead we can appeal to a corollary of  $\kappa \otimes \kappa = \kappa$ , which is already available [19, Section 5] in Isabelle/ZF.

$\llbracket \text{well\_ord}(A, r); \text{InfCard}(|A|) \rrbracket \implies A \times A \approx A$

Thus we have  $\omega \times \omega \approx \omega$ : there is a bijection, which is also an injection, between  $\omega \times \omega$  and  $\omega$ . However, although an injection exists, we have no means of naming a specific bijection. Therefore, we conduct the entire proof of the axiom of choice under the assumption that some injection exists. The final theorem is existential, which will allow the assumption to be discharged.

We declare a locale to express this new assumption, calling the injection *fn*. Recall that *nat* is Isabelle/ZF's name for the ordinal  $\omega$ .

```
locale Nat_Times_Nat =
  fixes fn
  assumes fn_inj: "fn ∈ inj(nat*nat, nat)"
```

Proving that  $\text{enum}(fn, p)$  defines an injection from formulae into the naturals requires a straightforward double induction over formulae.

$"(\lambda p \in \text{formula}. \text{enum}(fn, p)) \in \text{inj}(\text{formula}, \text{nat})"$

Using the enumeration as a measure function, we find that the set of formulae is well-ordered.

$"\text{well\_ord}(\text{formula}, \text{measure}(\text{formula}, \text{enum}(fn)))"$

The functions defined below all have an argument *f*, which should range over injections from  $\omega \times \omega$  into  $\omega$ . In proofs, this injection will always be *fn* from locale *Nat\_Times\_Nat*. The definiens of a constant definition cannot refer to *fn* because it is a variable.

### 10.3. Defining the well-ordering on $\text{DPow}(A)$

The set  $\text{DPow}(A)$  consists of those subsets of *A* that can be defined by a formula, possibly using elements of *A* as parameters (Section 6.5). We can define a well-ordering on  $\text{DPow}(A)$  from one on *A*. We get a well-ordering on formulae from their injection into the natural numbers. To handle the parameters, we define a well-ordering for environments – lists over *A* – and combine it with the well-ordering of formulae. A subset of *A* might be definable in more than one way; to make a unique choice, we map environment/formula pairs to ordinals. The well-ordering on environment/formula pairs is the lexicographic product (given by *rmult*) of the well-orderings on lists (*rlist*) and formulae (*measure*).

```
"env_form_r(f, r, A) ==
  rmult(list(A), rlist(A, r),
        formula, measure(formula, enum(f)))"
```

Using existing theorems, it is trivial to prove that this construction well-orders the set  $\text{list}(A) \times \text{formula}$ .

$"\text{well\_ord}(A, r) \implies \text{well\_ord}(\text{list}(A) \times \text{formula}, \text{env\_form\_r}(fn, r, A))"$

The order type of the resulting well-ordering yields a map (given by *ordermap*) from environment/formula pairs into the ordinals. For each member of  $\text{DPow}(A)$ , the minimum such ordinal will determine its place in the well-ordering.

```
"env_form_map(f, r, A, z) ==
  ordermap(list(A) × formula, env_form_r(f, r, A)) ` z"
```

If  $r$  well-orders  $A$ , and  $X$  is a definable subset of  $A$ , then let us define  $DPow\_ord(f, r, A, X, k)$  to hold if  $k$  corresponds to some definition of  $X$ ; that is, informally,  $k$  defines  $X$ .

```
"DPow_ord(f, r, A, X, k) ==
  ∃ env ∈ list(A). ∃ p ∈ formula.
    arity(p) ≤ succ(length(env)) &
    X = {x ∈ A. sats(A, p, Cons(x, env))} &
    env_form_map(f, r, A, (env, p)) = k"
```

Similarly, let us define  $DPow\_least(f, r, A, X)$  to be the smallest ordinal defining  $X$ .

```
"DPow_least(f, r, A, X) == μk. DPow_ord(f, r, A, X, k)"
```

Since  $k$  determines  $env$  and  $p$ , we find that an ordinal can define at most one element of  $DPow(A)$ .

```
"[DPow_ord(fn, r, A, X, k); DPow_ord(fn, r, A, Y, k); well_ord(A, r)]
  ⇒ X=Y"
```

We also find that every element of  $DPow(A)$  is defined by some ordinal, given by  $DPow\_least$ .

```
"[X ∈ DPow(A); well_ord(A, r)]
  ⇒ DPow_ord(fn, r, A, X, DPow_least(fn, r, A, X))"
```

Now  $DPow\_least$  can serve as a measure function to define the well-ordering on  $DPow(A)$ .

```
"DPow_r(f, r, A) == measure(DPow(A), DPow_least(f, r, A))"
```

Using general facts about relations defined by measure functions, we easily find that  $DPow(A)$  is well-ordered.

```
"well_ord(A, r) ⇒ well_ord(DPow(A), DPow_r(fn, r, A))"
```

#### 10.4. Well-ordering $L_\alpha$ in the limit case

The proof that  $L_\alpha$  is well-ordered appeals to transfinite induction on the ordinal  $\alpha$ . The induction hypothesis is that  $L_\xi$  is well-ordered if  $\xi < \alpha$ . In the limit case,  $L_\alpha = \bigcup_{\xi < \alpha} L_\xi$ . Recall (see Section 6.9.2) that the **L**-rank  $\rho(x)$  of  $x$  is the least  $\alpha$  such that  $x \in L_{\alpha+1}$ . If  $\alpha$  is a limit ordinal then we order elements of  $L_\alpha$  first by their **L**-ranks; if two elements have the same **L**-rank, say  $\xi$ , then we order them using the existing well-ordering of  $L_{\xi+1}$ .

In the Isabelle formalization,  $i$  is the limit ordinal and  $r(j)$  denotes the well-ordering of  $L_{set}(j)$ .

```
"rlimit(i, r) ==
  if Limit(i) then
    {z ∈ Lset(i) × Lset(i).
      ∃ x' x. z = (x', x) &
        (lrank(x') < lrank(x) |
         lrank(x') = lrank(x) & (x', x) ∈ r(succ(lrank(x))))}
  else 0"
```

We can prove that the limit ordering is linear, provided that the orderings of previous stages are also linear.

$$\begin{aligned} & \llbracket \text{Limit}(i); \forall j < i. \text{linear}(\text{Lset}(j), r(j)) \rrbracket \\ & \implies \text{linear}(\text{Lset}(i), rlimit(i, r)) \end{aligned}$$

Under analogous conditions, the  $rlimit(i, r)$  is a well-ordering of  $\text{Lset}(i)$ . The proofs are straightforward, and I have omitted many details.

$$\begin{aligned} & \llbracket \text{Limit}(i); \forall j < i. \text{well\_ord}(\text{Lset}(j), r(j)) \rrbracket \\ & \implies \text{well\_ord}(\text{Lset}(i), rlimit(i, r)) \end{aligned}$$

### 10.5. Transfinite definition of the well-ordering for $\mathbf{L}$

The well-ordering on  $\mathbf{L}$  is defined by transfinite recursion. The Isabelle definition refers to the cryptic *transrec* operator, so let us pass directly to the three immediate consequences of that definition. For the base case, the well-ordering is the empty relation.

$$\mid \text{"L\_r}(f, 0) = 0"$$

For the successor case, the well-ordering is given by applying  $DPow\_r$  to the previous level.

$$\mid \text{"L\_r}(f, succ(i)) = DPow\_r(f, L\_r(f, i), Lset(i))"$$

For the limit case, the well-ordering is given by  $rlimit$ .

$$\mid \text{"Limit}(i) \implies L\_r(f, i) = rlimit(i, L\_r(f))"$$

Thanks to the results proved above, a simple transfinite induction proves that  $L\_r(fn, i)$  well-orders the constructible level  $Lset(i)$ .

$$\text{"Ord}(i) \implies \text{well\_ord}(\text{Lset}(i), L\_r(fn, i))"$$

Note that this theorem refers to  $fn$ , an injection from  $\omega \times \omega$  into  $\omega$ . Recall (see Section 10.2) that we know such that such functions exist, but we have not defined a specific one. We have been able to prove our theorems by working in a locale that assumes the existence of  $fn$ . Now we can eliminate the assumption. We use an existential quantifier to hide the well-ordering in the previous theorem, so that  $fn$  no longer appears. Then, by the mere existence of such an injection, it follows that every  $Lset(i)$  can be well-ordered.

$$\text{"Ord}(i) \implies \exists r. \text{well\_ord}(\text{Lset}(i), r)"$$

To wrap things up, let us package the axiom  $\mathbf{V} = \mathbf{L}$  as a locale.

$$\begin{aligned} & \text{locale } V\_equals\_L = \\ & \quad \text{assumes } VL: \text{"L}(x) \end{aligned}$$

The axiom of choice – in the guise of the well-ordering theorem – is a trivial consequence of the previous results.

$$\text{theorem (in } V\_equals\_L) \text{ AC: } \llbracket \exists r. \text{well\_ord}(x, r) \rrbracket$$

## 11. Conclusions

What has been accomplished? I have mechanized the proof of the relative consistency of the axiom of choice, largely following a standard textbook presentation. The formal proof is much longer than the textbook version because it is complete in all details and uses no metatheoretical reasoning.

As noted in Section 2, Gödel's proof comprises four tasks, which we can now express more precisely:

1. defining the class  $\mathbf{L}$  within ZF;
2. proving, for every ZF axiom  $\phi$ , that  $\phi^{\mathbf{L}}$  is a ZF theorem;
3. proving  $(\mathbf{V} = \mathbf{L})^{\mathbf{L}}$  in ZF;
4. proving that  $\text{ZF} + \mathbf{V} = \mathbf{L}$  implies the axiom of choice.

The proof that  $\mathbf{L}$  satisfies  $\mathbf{V} = \mathbf{L}$  is by far the largest and most difficult part of the development. It involves proving  $\mathbf{L}$  to be absolute, which requires converting every concept used in its definition into relational form, and proving absoluteness. The sheer number of concepts is an obstacle, and some of them are hard to express in relational form, especially those involving recursion. Most of the relations have to be re-expressed using an internal datatype of formulae.

My formalization has two limitations. First, I am not able to prove that  $\mathbf{L}$  satisfies the axiom scheme of comprehension. Although Isabelle/ZF handles schematic proofs easily, the proof of comprehension for the formula  $\phi$  requires an instance of the reflection theorem for  $\phi$ . Each instance of comprehension therefore has a different proof, and must be proved separately. The reflection theorem is proved by induction (at the metalevel) on the structure of  $\phi$ ; thus, all these proofs are instances of one algorithm, and they are generated by nearly identical proof scripts [18]. The inability to prove the comprehension scheme makes the absoluteness proofs harder: every necessary instance of comprehension is listed. Instantiating these proofs to  $\mathbf{L}$  has required proving that each of those instances holds in  $\mathbf{L}$ . There are about thirty-five such instances.

My formalization has another limitation. The proof that  $\mathbf{L}$  satisfies  $\mathbf{V} = \mathbf{L}$  cannot be combined with the proof that  $\mathbf{V} = \mathbf{L}$  implies the axiom of choice in order to conclude that  $\mathbf{L}$  satisfies the axiom of choice. The reason is that the two instances of  $\mathbf{V} = \mathbf{L}$  are formalized differently: one is relativized, and the other is not. Here I have followed the textbook proofs, which prove that  $\mathbf{V} = \mathbf{L}$ , declare that the axiom of constructibility can be assumed, and proceed to derive the consequences of that axiom.

We could remedy both limitations by tackling the whole problem in a quite different way, by formalizing set theory as a proof system and working entirely in the metatheory. I leave this as a challenge for the theorem-proving community. A by-product of the work is a general theory of absoluteness for arbitrary class models of ZF. It could be used for other formal investigations of inner models. Future investigators might also try formalizing the proof that  $\mathbf{L}$  satisfies the generalized continuum hypothesis and the combinatorial principle.

*Acknowledgements.* Krzysztof Grąbczewski devoted much effort to an earlier, unsuccessful, attempt to formalize this material. Work on Isabelle is supported by the U.K.'s Engineering and Physical Sciences Research Council, grant GR/M75440. Markus Wenzel greatly improved Isabelle's *locale* construct to support these proofs. Kenneth Kunen gave advice that helped in my formalization of the reflection theorem, and the referee made a number of valuable comments on this paper.

Appendix A.

The appendix to the paper was automatically generated by Isabelle/ZF, and presents the full mechanical development. It is available at:

<http://www.lms.ac.uk/jcm/6/lms2003-001/appendix-a/>.

References

1. N. G. DE BRUIJN, ‘Lambda calculus notation with nameless dummies, a tool for automatic formula manipulation, with application to the Church–Rosser theorem’, *Indag. Math.* 34 (1972) 381–392. 215
2. B. A. DAVEY and H. A. PRIESTLEY, *Introduction to lattices and order* (Cambridge University Press, 1990). 230
3. S. FEFERMAN, *et al.*, eds, *Kurt Gödel: collected works*, vol. II. (Oxford University Press, 1990). 247
4. KURT GÖDEL, ‘The consistency of the axiom of choice and of the generalized continuum hypothesis’, [3] 26–27; first published in *Proc. Nat. Acad. Sci. USA* (1938) 556–557. 198
5. KURT GÖDEL, ‘The consistency of the axiom of choice and of the generalized continuum hypothesis with the axioms of set theory’, [3] 33–101; first published by Princeton University Press, 1940. 198, 199, 201, 215
6. KURT GÖDEL, ‘Consistency proof for the generalized continuum hypothesis’, [3] 27–32; first published in *Proc. Nat. Acad. Sci. USA* (1939) 220–224. 198, 215
7. PAUL R. HALMOS, *Naive set theory* (Van Nostrand, 1960). 201
8. FLORIAN KAMMÜLLER, MARKUS WENZEL and LAWRENCE C. PAULSON ‘Locales: a sectioning concept for Isabelle’, *Theorem proving in higher order logics: TPHOLS ’99*, Lecture Notes in Comput. Sci. 1690 (ed. Yves Bertot, Gilles Dowek, André Hirschowitz, Christine Paulin and Laurent Théry, Springer, 1999) 149–165. 208
9. KENNETH KUNEN, *Set theory: an introduction to independence proofs* (North-Holland, 1980). 200, 201, 203, 204, 211, 212, 219, 220, 221, 222, 227, 241
10. E. MENDELSON, *Introduction to mathematical logic*, 4th edn (Chapman and Hall, 1997). 200
11. TOBIAS NIPKOW, LAWRENCE C. PAULSON and MARKUS WENZEL, *Isabelle/HOL: a proof assistant for higher-order logic*, Lecture Notes in Comput. Sci. Tutorial 2283 (Springer, 2002). 204
12. LAWRENCE C. PAULSON, ‘The foundation of a generic theorem prover’, *J. Automat. Reasoning* 5 (1989) 363–397. 223
13. LAWRENCE C. PAULSON, ‘Set theory for verification: I. From foundations to functions’, *J. Automat. Reasoning* 11 (1993) 353–389. 204, 211
14. LAWRENCE C. PAULSON, *Isabelle: a generic theorem prover*, Lecture Notes in Comput. Sci. 828 (Springer, 1994). 204
15. LAWRENCE C. PAULSON, ‘Set theory for verification: II. Induction and recursion’ *J. Automat. Reasoning* 15 (1995) 167–215. 212, 213, 214, 230, 232, 233, 234, 235

16. LAWRENCE C. PAULSON, ‘Proving properties of security protocols by induction’, *10th Computer Security Foundations Workshop* (IEEE Computer Society Press, 1997) 70–83. 204
17. LAWRENCE C. PAULSON, ‘A fixedpoint approach to (co)inductive and (co)datatype definitions’, *Proof, language, and interaction: essays in honor of Robin Milner* (ed. Gordon Plotkin, Colin Stirling, and Mads Tofte, MIT Press, 2000) 187–211. 204, 230, 231
18. LAWRENCE C. PAULSON, ‘The reflection theorem: a study in meta-theoretic reasoning’, [22] 377–391. 198, 222, 223, 224, 246
19. LAWRENCE C. PAULSON and KRZYSZTOF GRĄBCZEWSKI, ‘Mechanizing set theory: cardinal arithmetic and the axiom of choice’, *J. Automat. Reasoning* 17 (1996) 291–323. 204, 243
20. DAG PRAWITZ, ‘Ideas and results in proof theory’, *Second Scandinavian Logic Symposium* (ed. J. E. Fenstad, North-Holland, 1971) 235–308. 202
21. MARTIN STRECKER, ‘Formal verification of a java compiler in Isabelle, [22] 63–77. 204
22. ANDREI VORONKOV, ed. *Automated deduction – CADE-18 International Conference*, Lecture Notes in Artificial Intelligence 2392 (Springer, 2002). 248

Lawrence C. Paulson   lcp@cl.cam.ac.uk  
<http://www.cl.cam.ac.uk/users/lcp>

Computer Laboratory  
University of Cambridge  
J J Thomson Avenue  
Cambridge CB3 0FD  
England