

# Proof-producing reflection for HOL with an application to model polymorphism

Benja Fallenstein<sup>1</sup> and Ramana Kumar<sup>2</sup>

<sup>1</sup> Machine Intelligence Research Institute

<sup>2</sup> Computer Laboratory, University of Cambridge

**Abstract.** We present a reflection principle of the form “If  $\ulcorner\varphi\urcorner$  is provable, then  $\varphi$ ” implemented in the HOL4 theorem prover, assuming the existence of a large cardinal. We use the large-cardinal assumption to construct a model of HOL within HOL, and show how to ensure  $\varphi$  has the same meaning both inside and outside of this model. Soundness of HOL implies that if  $\ulcorner\varphi\urcorner$  is provable, then it is true in this model, and hence  $\varphi$  holds. We additionally show how this reflection principle can be extended, assuming an infinite hierarchy of large cardinals, to implement *model polymorphism*, a technique designed for verifying systems with self-replacement functionality.

## 1 Introduction

*Reflection principles* of the form<sup>3</sup> “if  $\ulcorner\varphi\urcorner$  is provable, then  $\varphi$ ” have long been of interest in logic and set theory (see, e.g., Franzén [5] and Jech [13]). In this paper, we show how to implement a reflection principle for HOL in the HOL4 theorem prover [21], using a novel approach for establishing a correspondence between the logic and an internal model. Such a reflection principle does not come for free, since by Gödel’s second incompleteness theorem, HOL cannot prove itself consistent [6]. We pay with an assumption about the existence of a “large enough” HOL type. But we endeavour to ensure that this assumption has the same content as the assumption, commonly studied in set theory [13], of the existence of a strongly inaccessible cardinal.

Reflection is about trying to fit a logic inside itself, so one has to keep two instances of the logic in mind separately. We use the term *inner HOL* to refer to the object logic (the HOL that is formalised) and *outer HOL* for the meta logic (the HOL in which inner HOL is formalised). We build upon Harrison’s formalisation of HOL in itself [11] that was extended in our previous work [15, 16] to support defined constants. The first reflection principle we implement uses a model of inner HOL provided by a large-cardinal assumption, together with the previously-proved soundness theorem, to show that provability of a proposition in inner HOL implies its truth in outer HOL.

---

<sup>3</sup>  $\ulcorner\varphi\urcorner$  refers to  $\varphi$  as a syntactic object, represented e.g. by its Gödel number or by an abstract syntax tree.

This kind of reflection principle is asymmetric: the large-cardinal assumption used in outer HOL to justify reflection cannot be used again in inner HOL to justify further reflection. One setting in which it is useful to lessen this asymmetry is in constructing and verifying a system that can replace itself (including the replacement mechanism) by a new version while nevertheless always satisfying some safety property. We have extended the reflection principle above to an implementation of *model polymorphism* [3], which enables the use of reflective reasoning multiple times to verify such self-modifying systems.

Our specific contributions in this paper are as follows:

- A simple approach to defining, in outer HOL, a partial embedding of outer HOL in inner HOL (§3). We re-use a single polymorphic constant, `to_inner`, to embody the embedding at (almost) all outer HOL types.
- An algorithm for producing theorems that relate the semantics of inner HOL terms to their outer HOL counterparts via the embedding (§4). The method applies to terms that include user-defined constants and types, and therefore supports construction of a semantic interpretation of constants based on replaying outer definitions in inner HOL (§5). Combined with the soundness theorem for inner HOL, this gives us a reflection principle for HOL (§7).
- An extension of the reflection principle above to support model polymorphism [3] (§8).
- A reduction of the unavoidable assumption on the proof (in previous work) of soundness and consistency of HOL within HOL to a traditionally-stated large-cardinal assumption (§6).

All the work we present has been implemented in the HOL4 theorem prover, and is available online at <https://github.com/CakeML/hol-reflection>.

*Synopsis* Our reflection principle has two legs: (1) a soundness theorem for HOL (from previous work) asserting that a provable formula is true in all models, and (2) a formal correspondence between objects in a particular, reflective model (constructed assuming a large cardinal) and their counterparts outside. The reflection principle works as follows: if a formula  $\ulcorner \varphi \urcorner$  is provable, the soundness theorem asserts that  $\ulcorner \varphi \urcorner$  is true in the reflective model of HOL, and via the correspondence we may conclude that  $\varphi$  is true outside the model. The bulk of our work in this paper is in building the reflective model and establishing the correspondence, particularly in supporting user-defined constants.

## 2 Background: inner HOL

The starting point for studying reflection principles is a formalisation of a logic inside itself. We use higher-order logic (HOL), the logic of the HOL Light theorem prover [12]. Our formalisation of HOL within itself is described carefully in previous work [15, 16], but we summarise it briefly in this section.

HOL has the syntax of the simply-typed lambda calculus<sup>4</sup> and a well-known semantics in Zermelo set theory. The judgements of the logic are sequents,  $(thy, hs) \vdash c$ , where  $c$  is a formula (a term of Boolean type),  $hs$  is a set of formulas, and  $thy$  specifies the current set of axiom formulas and the current signature of constants and type operators. The semantics of each type is a non-empty set; the Boolean type is specified as a distinguished two-element set (`Boolset`, with elements `True` and `False`) and function types are specified as function spaces (a set of pairs forming a functional relation). The semantics of each term is an element of the semantics of its type. The semantics of each sequent is true, written  $(thy, hs) \models c$ , if the semantics of  $c$  is true whenever the semantics of all  $hs$  are true, in any model (an interpretation of constants satisfying the axioms) of  $thy$ .

Since the semantic objects are sets, we need a model of set theory to define the semantics of inner HOL. Our semantics does not pin down a particular model of set theory. Instead, it is polymorphic: we use an outer HOL type variable, usually  $\mu$ , for the universe of sets, and encode the axioms of set theory as a pair of assumptions, `is_set_theory mem` and  $\exists inf. is\_infinite mem inf$ , about a membership relation,  $(mem : \mu \rightarrow \mu \rightarrow bool)$ , on sets. All our semantic functions take the membership relation  $mem$  (and with it the universe  $\mu$ ) as a parameter, and most theorems about the semantics assume the set-theory axioms hold. However, for brevity we usually hide<sup>5</sup> these parameters and assumptions when they are obvious in context. We write  $x < y$  for infix application of the  $mem$  relation (i.e.,  $mem x y$ ).

The important semantic function in this paper is the semantics of terms, that is, the function that assigns an element of  $\mu$  to each term of inner HOL. We illustrate how it works with an example, considering the semantics of the inner HOL version of the term  $\lambda x. foo\ y$ , where the constant  $(foo : \alpha \rightarrow \beta)$  is instantiated at the types `bool` and `ind list`. The four parameters governing the semantics are: the membership relation  $mem$  (hidden), the signature ( $s$ ) of constants and type operators in the theory, the interpretation ( $i$ ) of constants and type operators, and the valuation ( $v$ ) of free term and type variables. We denote the valuation of type variables by  $v_{ty}$  whereas for term variables we use  $v_{tm}$ , and similar subscripts apply to the other semantic parameters. Our use of interpretations and valuations is intended to be conventional and unsurprising; we show the example here mainly to make our notation clear.<sup>6</sup>

```
termsem s_tm i v
  (Abs (Var "x" Bool)
    (Comb (Const "foo" (Fun Bool (Tyapp "list" [Tyapp "ind" []])))
      (Var "y" Bool))) =
  Abstract Boolset sil (\ sx. sfoo ' sy)
```

<sup>4</sup> Including, as in Gordon [9], polymorphism (type variables and polymorphic constants) and defined constants and type operators.

<sup>5</sup> This is akin to working within an Isabelle [23] locale which fixes  $\mu$  and  $mem$  and assumes `is_set_theory mem`. (We assume infinity only when necessary.)

<sup>6</sup> The abstract syntax here is the inner HOL rendition of our example term. `Bool` and `Fun a b` are abbreviations for `Tyapp "bool" []` and `Tyapp "fun" [a; b]`.

The semantics of the given lambda abstraction is a set-theoretic function (created with `Abstract`) – in particular the set of two pairs, `True`  $\mapsto$  `sfoo ' sy` and `False`  $\mapsto$  `sfoo ' sy`. The following equations show how we obtain the semantics of the applied type operator, the instantiated constant, and the free variable:

$$\begin{aligned} sil &= i_{ty} \text{"list"} [i_{ty} \text{"ind"} []] \\ sfoo &= i_{tm} \text{"foo"} [\text{Boolset}; sil] \\ sy &= v_{tm} (\text{"y"}, \text{Bool}) \end{aligned}$$

Inner HOL supports defined constants and type operators via the *current theory*, which is attached to each sequent. The inference system uses a concrete implementation of each theory as a list of *updates*. Such a list is called a *context* and every context can be viewed abstractly as a theory (written: `thyof ctxt`). Contexts are made from five kinds of update: new type operator, new constant, new axiom, new defined type operator, and new defined constant. The updates for defined type operators and constants have preconditions that require sequents proved in the previous context. If an update *upd* satisfies all the preconditions to be a valid update of the previous context *ctxt*, we say *upd updates ctxt*.

Both the within-theory inference rules and the theory-extending rules for making updates (except new axiom) have been proved sound with respect to the semantics. For the inference rules, the soundness theorem states<sup>7</sup> that every provable sequent is true:

$$\vdash (thy, hs) \vdash c \Rightarrow (thy, hs) \models c$$

For the extension rules, the soundness theorem states that there exists an extended interpretation of the new theory that continues to be a model of the theory. Both theorems require the `is_set_theory mem` assumption.

### 3 An inner copy of outer HOL

Given an outer HOL term, it is straightforward to write a function in ML—our function is called `term_to_deep`—that walks the structure of the term and builds the corresponding inner HOL term. For example, `term_to_deep`

turns `Suc x` into `Comb (Const "Suc" (Fun Num Num)) (Var "x" Num)`.

This syntactic connection is straightforward. But what is the relationship between the semantics of an outer HOL term and its inner counterpart? Let  $\ulcorner tm \urcorner$  stand for `term_to_deep(tm)`. Ideally, we would like a connection between `Suc x` and `termsem stm i v`  $\ulcorner$ `Suc x` $\urcorner$ . Such a connection would mean that the structure of outer HOL terms and types is replicated (indeed reflected) within the type

<sup>7</sup> To understand the three turnstiles:  $\vdash$  denotes provability in outer HOL (and applies to the whole formula),  $\vdash$  denotes provability in inner HOL, and  $\models$  states that a sequent is valid according to the semantics of inner HOL.

denoted by  $\mu$ . We cannot expect to reflect everything: in particular, we cannot reflect  $\mu$  within itself, nor anything depending on  $\mu$ . But we can cover all other outer HOL types and terms generically.

Our solution is to define a polymorphic constant, `to_inner`, which sends an outer HOL term to the element of  $\mu$  to which it is supposed to correspond. We show how to prove theorems of the form  $\text{termsem } s_{\text{tm}} \ i \ v \ \ulcorner \text{tm} \urcorner = \text{to\_inner } \text{ty } \text{tm}$ . The type of `to_inner ty` is  $\alpha \rightarrow \mu$ .<sup>8</sup> (We explain the *ty* argument shortly.)

We do not want to have to specify exactly how `to_inner` creates its (partial) copy of outer HOL. What is important is that the copy is faithful, which means `to_inner` is injective. We formalise this injectivity property precisely by saying that `to_inner` at type  $\alpha$  should be a bijection between  $\mathcal{U}(:\alpha)$  (everything of type  $\alpha$ ) and the set<sup>9</sup> of elements of some set  $x$  in  $\mu$ . Formally, we define a well-formedness condition:

$$\text{wf\_to\_inner } f \iff \exists x. \text{BIJ } f \ \mathcal{U}(:\alpha) \ \{ a \mid a \triangleleft x \}$$

Then we define `to_inner` as an arbitrary well-formed injection, using Hilbert choice<sup>10</sup> as follows:

$$\text{to\_inner } \text{ty} = \text{tag } \text{ty} \circ \varepsilon f. \text{wf\_to\_inner } f$$

The `tag ty` part of the definition wraps the set produced by the well-formed injection with a tag for the given inner HOL type, *ty*, thereby avoiding the possibility of inadvertently sending outer HOL terms with different types to the same element in  $\mu$ . The need for this tagging is explained in Section 5.

Since Hilbert choice only picks a well-formed injection *if one exists*, the usefulness of `to_inner` at any particular outer HOL type depends on our assuming (or being able to prove) `wf\_to\_inner (to_inner ty)`. The automation we describe in Section 4 produces theorems that assume these well-formedness conditions, and the automation in Section 5 proves almost all of them.

Of course, since the well-formedness condition on `to_inner` says it is not just an injection but a bijection, we can also go in the other direction, from inner HOL terms to their outer counterparts. Given a well-formed injection *ina*, we denote the set of terms of type  $\mu$  that are in its range—that is, the inner representation of the domain of *ina*—by `range ina`. The faithfulness of the representation is summarised in the following two theorems exhibiting invertibility.

$$\begin{aligned} \vdash \text{wf\_to\_inner } \text{ina} &\Rightarrow \forall x. \text{ina}^{-1} (\text{ina } x) = x \\ \vdash \text{wf\_to\_inner } \text{ina} &\Rightarrow \forall x. x \triangleleft \text{range } \text{ina} \Rightarrow \text{ina} (\text{ina}^{-1} x) = x \end{aligned}$$

<sup>8</sup> To be pedantic, `to_inner` also depends on the pervasive *mem* relation of the set theory.

<sup>9</sup> We follow the convention of treating predicates in outer HOL as sets. To be clear  $\mathcal{U}(:\alpha)$  is a term of type  $\alpha \rightarrow \text{bool}$ , and  $\{ a \mid a \triangleleft x \}$  is a term of type  $\mu \rightarrow \text{bool}$ . These sets-as-predicates are distinct from the Zermelo sets, i.e., the terms of type  $\mu$ .

<sup>10</sup> Also known as indefinite choice, the Hilbert choice principle in HOL provides a constant ( $\varepsilon$ ), usually written as a binder, together with the axiom  $(\exists x. P x) \Rightarrow P (\varepsilon x. P x)$  which holds for any predicate  $P$ .

Usually, but not always,  $ina$  will be  $to\_inner\ ty$  for some  $ty$ . We use  $to\_inner$  to reflect outer HOL at all types except for those that depend on  $\mu$  and except for the two primitive types of HOL: Booleans and function types. The primitive types must be treated differently because the semantics of HOL requires them to be interpreted by a distinguished two-element set and by set-theoretic function spaces, respectively. We provide specialised constants that map these types to their intended interpretations instead of an arbitrary set:

```
bool_to_inner b = if b then True else False
fun_to_inner ina inb f =
  Abstract (range ina) (range inb) ( $\lambda x. inb (f (ina^{-1} x))$ )
```

We have now seen how we intend to reflect outer HOL terms into inner HOL, using  $to\_inner\ ty$  injections. We next describe an algorithm that produces a certificate theorem for (almost) any HOL term, asserting that the semantics of the inner version of the term matches the reflection of the outer version. We develop this algorithm in two stages: first (next section), we ignore the interpretation of constants and the valuation of variables in the semantics and simply make assumptions about them; then (Section 5), we build an interpretation and valuation that satisfies these assumptions.

## 4 Proof-producing reflection

Let us begin with an example of the kind of theorem produced by the first stage of our automation. Given as input the term  $Suc\ x$ , we produce the following theorem that looks very long but whose components we will explain one by one. We call each such theorem a *certificate theorem* for the input term.

```
⊢ good_context mem s i ∧ wf_to_inner (to_inner Num) ∧
  lookup s_tm "Suc" = Some (Fun Num Num) ∧ lookup s_ty "num" = Some 0 ∧
  i_ty "num" [] = range (to_inner Num) ∧
  i_tm "Suc" [] = fun_to_inner (to_inner Num) (to_inner Num) Suc ∧
  v_tm ("x", Num) = to_inner Num x ⇒
  termsem s_tm i v ⌈Suc x⌋ = to_inner Num (Suc x)
```

First, look at the conclusion<sup>11</sup> (the last line): we have produced an equality between the semantics of the inner version of our input term and its reflection created with  $to\_inner$ . Under what assumptions? The assumptions come in five categories:

1. The assumption stated using `good_context`, which represents the pervasive `is_set_theory mem` assumption as well as some basic well-formedness conditions on the signature and interpretation.

<sup>11</sup> `Num` is an abbreviation for `Tyapp "num" []`.

2. `wf_to_inner` assumptions on all *base types* appearing in the input term (in this case, just `Num`). A base type is any type variable or any application of a non-primitive type operator. (For the primitives, Booleans and functions, we prove `wf_to_inner` once and for all.)
3. Signature (*s*) assumptions stating that all non-primitive constants and type operators in the input term (in this case `Suc` and *num*) have the same type/arity in inner HOL as in outer HOL. (The only primitive constant, equality, is assumed to have the correct type as part of `good_context`.)
4. Interpretation (*i*) assumptions stating that the inner versions of all the base types and constants of the input term are mapped to their reflections.
5. Valuation (*v*) assumptions stating that the inner versions of all type and term variables in the input term are mapped to their reflections.

The algorithm for producing such a theorem works by recursively traversing the structure of the input term. We build a theorem like the one above, equating the semantics of an inner term to the reflection of its outer counterpart, for each subterm in the input term starting at the leaves and progressing bottom-up. When we encounter non-primitive type operators and constants we add signature and interpretation assumptions as required, and similarly add valuation assumptions for free variables. The substitution of type variables used to instantiate a polymorphic constant is easily reflected from outside to inside.

This algorithm works because the semantics itself is recursive. For example, shown below is a theorem about the semantics that our algorithm uses when it encounters a combination (function application) term. Notice that the theorem makes two assumptions of the same form as its conclusion – these correspond to recursive calls in the algorithm.

$$\begin{aligned} \vdash \text{termsem } s_{\text{tm}} \ i \ v \ f\text{tm} = \text{fun\_to\_inner } \text{ina} \ \text{inb} \ f \wedge \text{termsem } s_{\text{tm}} \ i \ v \ x\text{tm} = \text{ina} \ x \wedge \\ \text{wf\_to\_inner } \text{ina} \wedge \text{wf\_to\_inner } \text{inb} \Rightarrow \\ \text{termsem } s_{\text{tm}} \ i \ v \ (\text{Comb } f\text{tm} \ x\text{tm}) = \text{inb} \ (f \ x) \end{aligned}$$

The analogous theorem for lambda abstractions requires us to prove connections between inner HOL types and their reflections (via `range`). Therefore, we have a similar recursive algorithm for types as the one described above for terms.

So far, our certificate theorems leave the semantic parameters (*s*, *i*, and *v*) as free variables but make various assumptions about them. Our aim in the next section is to show how we can instantiate these parameters in such a way that most of the assumptions become provable. In other words, we show how to build a reflective model that satisfies the assumptions in each of the categories above.

## 5 Building a reflective interpretation and valuation

*Signature assumptions* In outer HOL, types like *num* and constants like `Suc` are all user-defined. Logically speaking, there is a context (typically not represented explicitly in the theorem prover) that lists the sequence of updates made to produce the current environment of defined constants. An appropriate signature (*s*)

for a certificate theorem is one that corresponds to some context for the input term, which will naturally satisfy all the signature assumptions. We require the user of our automation to build an explicit representation of the desired context (although we provide tools to help with this), and with that information proceed to construct an interpretation and valuation to satisfy the assumptions of the certificate theorem.

*Algorithm for building the model* The idea is to reflect the updates from the outer context into the inner inference system, creating an inner context, and then, update by update, build an interpretation of the inner context. The interpretation we build must be reflective, which means it must satisfy the certificate theorem’s assumptions, namely: all the defined constants and types appearing in the input term are mapped to the reflections of their outer versions. We build the reflective interpretation by: (1) starting with the interpretation asserted to exist by the soundness theorem, and (2) making a finite number of modifications for each update in the context. We work recursively up the context applying the relevant modifications for the last-added update at each stage.

*Interpretation assumptions* Each update introduces some number of type operators, constants, and axioms<sup>12</sup> to the theory. Each update thereby induces some *constraints* on an interpretation for the interpretation to be reflective. For example, the update that defines  $\ulcorner \text{Suc} \urcorner$  has an associated constraint that this constant is interpreted as `to_inner (Fun Num Num) Suc`.

For polymorphic constants and type operators with arguments, the constraints are more involved. In general, a constraint is induced by each instance of the update and constrains the corresponding instance of the introduced constants. Importantly, we only consider instances of the update that correspond to the finitely many instances of the defined constants in the input term. We require `to_inner` to produce distinct reflections of distinct types, because otherwise we cannot reliably distinguish different instances of an update.

To show that an interpretation continues to be a model after being constrained, we require that at each constrained instance the axioms of the update are satisfied. We can prove this for the constraints we build (which map things to their reflections) because the axioms of the update being replayed are true in outer HOL. For each constrained instance of an axiom, we use the algorithm from the previous section to generate a certificate theorem with assumptions that are all easy to discharge because the constraint is reflective.

*wf\_to\_inner assumptions* To prove the `wf_to_inner` assumptions (for defined type operators) on our certificate theorems, we use the fact that in outer HOL every defined type is represented by a predicate on a previously defined (or primitive) type. The same is true in inner HOL as we build up our interpretation, so the

<sup>12</sup> When a new constant or type operator is defined, theorems produced by the definition are considered axioms of the resulting theory. The other source of axioms is new-axiom updates, which we do not support since they are not sound in general.



`wf_to_inner` assumption on a defined type reduces to the `wf_to_inner` assumption on its representing type. These assumptions propagate back recursively to the base case. The only `wf_to_inner` assumptions left after this process are for type variables in the input term, and for the type of individuals (*ind*) – this last assumption may be introduced if not present.

*Base case of the algorithm* To finish describing how we build a constrained model that satisfies the assumptions of the certificate theorem, only the base case of the recursive algorithm remains. What model do we start with before replaying all the updates in the context? In our previous work [15, 16] on the soundness and consistency of HOL, we showed that there is a model for the base context of HOL (assuming the set-theory axioms, including infinity). We use this model in the base case of the algorithm for building a reflective interpretation, modulo some subtleties concerning the Hilbert choice operator.

The base context, `hol_ctxt`, includes the primitive types and constants (Booleans, functions, equality) and the three axioms of HOL (extensionality, choice, and infinity). To state the axioms, the base context also includes some defined constants (conjunction, implication, etc.), and, importantly, the Hilbert choice constant, and the type of individuals. Although most constraints are dealt with before the base case, constraints on Hilbert choice remain. Therefore, we extend the proof that `hol_ctxt` has a model to show that in fact there is a model satisfying any finite number of constraints on the interpretation of Hilbert choice. We also modify the model to use `range (to_inner lnd)` for the type of individuals.

*Valuation assumptions* So far we have described the algorithm for building an interpretation (*i*) to satisfy the interpretation assumptions on certificate theorems. To build a valuation (*v*) satisfying the valuation assumptions, we follow a similar approach based on constraining an arbitrary valuation so that it maps the free variables that occur in the input term to their reflections. We have not covered all the gory details of these algorithms, but hope to have shared the important insights. The ML code for our automation is around 1800 lines long, supported by around 1000 lines of proof script about constrained interpretations.

## 6 Set theory from a large-cardinal assumption

Up until now, we have been working under the assumption `is_set_theory mem`, as used in previous work [15, 16], which consists of the set-theoretic axioms of extensionality, specification, pairing, union, and power sets. To produce models of `hol_ctxt` and all its extensions, we also require the set-theoretic axiom of infinity (which is implied by `wf_to_inner (to_inner lnd)`). One of our contributions is to clarify that these assumptions are implied by a more traditionally-stated large-cardinal assumption.

Our large-cardinal assumption is strictly stronger than what is necessary for a model of HOL. We leave proving an equivalence between `is_set_theory mem` plus infinity and a large-cardinal assumption for future work, and here just prove

implication from a strongly inaccessible cardinal. The point is to build confidence in our specification of set theory, by showing it is implied by a traditional specification (of a strong inaccessible).

We make use of Norrish and Huffman’s formalisation [20] of cardinals and ordinals in HOL4. We write  $s \prec t$  to mean the predicate  $s$  has smaller cardinality than the predicate  $t$ , i.e., there is an injection from  $s$  to  $t$  but not vice versa. Remember that sets and predicates are usually identified in HOL; we use the usual notation for sets like  $x \in s$  ( $x$  is an element of  $s$ ) and  $f \ " \ s$  (the image of  $s$  under  $f$ ) when  $s$  is a predicate.

The reduction theorem we have proved is that if some (outer HOL) type, say  $\mu$ , is a *strongly inaccessible cardinal*, then there is a membership relation  $mem$  on that type that satisfies `is_set_theory` and `is_infinite`. Furthermore, according to this  $mem$ , for every predicate  $s$  on the type  $\mu$  which is strictly smaller than  $\mathcal{U}(:\mu)$ , there is a Zermelo set  $x$  whose elements are exactly the extension of  $s$ . In this section, we are explicit with all uses of the type variable  $\mu$  and the variable  $mem$ . The formal statement of our reduction theorem is:

$$\begin{aligned} \vdash \text{strongly\_inaccessible } \mathcal{U}(:\mu) \Rightarrow \\ \exists mem. \\ \text{is\_set\_theory } mem \wedge \\ (\forall s. s \prec \mathcal{U}(:\mu) \Rightarrow \exists x. s = \{ a \mid a \prec x \}) \wedge \\ \exists inf. \text{is\_infinite } mem \text{ inf} \end{aligned}$$

It is well known that the existence of a strongly inaccessible cardinal gives rise to a model of set theory [13, Lemma 12.13], and the proof of our reduction theorem contains no surprises. (The main lemma is that a strongly inaccessible cardinal is in bijection with its smaller subsets.) Our focus here is on the definition of `strongly_inaccessible`, aiming to accurately capture traditional usage.

A cardinal is called “strongly inaccessible” if it is uncountable, a regular cardinal, and a strong limit cardinal [13]. Assuming the axiom of choice (which we do as we are working in HOL), a cardinal  $X$  is regular iff it cannot be expressed as the union of a set smaller than  $X$  all of whose elements are also smaller than  $X$  [13, Lemma 3.10]. Formally:

$$\begin{aligned} \text{regular\_cardinal } X \iff \\ \forall x f. \\ x \subseteq X \wedge x \prec X \wedge (\forall a. a \in x \Rightarrow f \ a \subseteq X \wedge f \ a \prec X) \Rightarrow \\ \bigcup (f \ " \ x) \prec X \end{aligned}$$

A cardinal is a strong limit if it is larger than the power set of any smaller cardinal. This is straightforward to formalise:

$$\text{strong\_limit\_cardinal } X \iff \forall x. x \subseteq X \wedge x \prec X \Rightarrow \mathcal{P}(x) \prec X$$

A cardinal is countable if it can be injected into the natural numbers (this is already defined in HOL4’s standard library). With these three ideas formalised, we define `strongly_inaccessible` as follows:

$$\begin{aligned} \text{strongly\_inaccessible } X \iff \\ \text{regular\_cardinal } X \wedge \text{strong\_limit\_cardinal } X \wedge \neg \text{countable } X \end{aligned}$$

## 7 Proving reflection principles

We now have enough machinery in place to exhibit a reflection principle for HOL. In the examples that follow, we use a two-place predicate `Safe t v` to construct our input propositions, to match the example in Section 8. However, `Safe` can be considered as a placeholder for any two-place predicate. As a concrete, simple example, set `Safe t v  $\iff$  v  $\neq$  Suc t`.

First, let us see how the reflection principle works on the input proposition  $\forall t. \text{Safe } t \ 0$ . Combining the automation from Sections 4 and 5, we can construct a certificate theorem with almost all of the assumptions proved:

$$\begin{aligned} &\vdash \text{is\_set\_theory } mem \wedge \text{wf\_to\_inner (to\_inner Ind)} \Rightarrow \\ &\quad \text{termsem (sigof inner\_ctxt)\_{tm} constrained\_model constrained\_valuation} \\ &\quad \ulcorner \forall t. \text{Safe } t \ 0 \urcorner = \\ &\quad \text{bool\_to\_inner } (\forall t. \text{Safe } t \ 0) \end{aligned}$$

Here `constrained_model` and `constrained_valuation` are built as described in Section 5: they are a reflective model for the context (`inner_ctxt`) that defines  $\ulcorner \text{Safe} \urcorner$ , and a reflective valuation<sup>13</sup> for the input proposition.

Next, consider the inner HOL sequent corresponding to provability of our input proposition:

$$(\text{thyof inner\_ctxt, []}) \vdash \ulcorner \forall t. \text{Safe } t \ 0 \urcorner$$

The soundness theorem for inner HOL states that if this sequent is provable, then the semantics of its conclusion is true. Our certificate theorem above already tells us that the semantics of the conclusion is equal to the reflection of the input proposition. By definition of `bool_to_inner`, if the reflection of a proposition is true then that proposition holds, hence we obtain:

$$\begin{aligned} &\vdash \text{is\_set\_theory } mem \wedge \text{wf\_to\_inner (to\_inner Ind)} \Rightarrow \\ &\quad (\text{thyof inner\_ctxt, []}) \vdash \ulcorner \forall t. \text{Safe } t \ 0 \urcorner \Rightarrow \forall t. \text{Safe } t \ 0 \end{aligned}$$

There would be additional `wf_to_inner` assumptions for each additional type variable in the input term.

The final step is to replace the `is_set_theory mem` assumption with the large-cardinal assumption that provides us with a model of set theory. In this model, we can replace the `wf_to_inner` assumption on `ind` by a lower bound on the cardinality of  $\mu$ . Similarly, any `wf_to_inner` assumption on a type variable, say  $\alpha$ , could be replaced by an assumption of the form  $\mathcal{U}(:\alpha) \prec \mathcal{U}(:\mu)$ . Such assumptions can be satisfied by large enough  $\mu$ , except in the case  $\alpha = \mu$ ; hence, it is important that  $\mu$  not occur in the input proposition. The resulting theorem no longer contains any occurrences of `mem`, not even hidden ones.

$$\begin{aligned} &\vdash \text{strongly\_inaccessible } \mathcal{U}(:\mu) \wedge \mathcal{U}(:ind) \prec \mathcal{U}(:\mu) \Rightarrow \\ &\quad (\text{thyof inner\_ctxt, []}) \vdash \ulcorner \forall t. \text{Safe } t \ 0 \urcorner \Rightarrow \forall t. \text{Safe } t \ 0 \end{aligned}$$

<sup>13</sup> Any valuation would do for this input since it has no free type or term variables.

Thus we have shown that provability of our input proposition implies its truth, assuming the existence of a large cardinal that is larger than the type of individuals in outer HOL. We can prove a theorem like this for any input proposition (that does not mention  $\mu$ ), including less obvious and even false propositions.

The reflection principle above can be generalised to a *uniform reflection principle* [4], which, for an input proposition with a free natural-number variable, requires an inner-HOL proof about only the relevant value of the variable. We write  $\ulcorner \bar{v} \urcorner$  for the inner HOL numeral corresponding to the value of the outer HOL variable  $v$ ; for example, if  $v = 1$  in outer HOL, then  $\ulcorner \text{Safe } t \bar{v} \urcorner$  denotes the term  $\ulcorner \text{Safe } t (\text{Suc } 0) \urcorner$ . The uniform reflection principle for the predicate  $\forall t. \text{Safe } t v$  is:

$$\begin{aligned} &\vdash \text{strongly\_inaccessible } \mathcal{U}(:\mu) \wedge \mathcal{U}(:ind) \prec \mathcal{U}(:\mu) \Rightarrow \\ &\quad \forall v. (\text{thyof inner\_ctxt}, []) \vdash \ulcorner \forall t. \text{Safe } t \bar{v} \urcorner \Rightarrow \forall t. \text{Safe } t v \end{aligned}$$

Since  $v$  is quantified in outer HOL, this theorem encapsulates infinitely many reflection theorems of the previous kind. To implement uniform reflection, we define an outer-HOL function `quote : num  $\rightarrow$  term` which can be spliced into the result of `term_to_deep` to provide terms of the form  $\ulcorner \bar{v} \urcorner$ . It is straightforward to show that the semantics of  $\ulcorner \bar{v} \urcorner$  in a reflective valuation is equal to the reflection of  $v$ , and uniform reflection follows.

## 8 An implementation of model polymorphism

One application of reflection principles is in designing and verifying systems that include mechanisms for self-replacement. For concreteness, consider an operating system intended to satisfy a certain safety property (e.g., a certain file is never overwritten), but also with a mechanism for replacing itself by an arbitrary updated version. For this system to be safe, the replacement mechanism must be restricted to prevent replacement by an unsafe update. A simple restriction would be to require a proof that the updated version is safe until replacement is invoked, together with a syntactic check that the updated version's replacement mechanism (including the proof checker) is unchanged (or that the replacement mechanism is removed altogether). To verify this system, we need to know that the system's proof checker is sound (only admits valid proofs) and we need a reflection principle for its logic (the conclusions of valid proofs are true). But we only need to establish these properties once.

Things get more interesting if we want to allow updates that might change the replacement mechanism. A more general replacement mechanism simply requires a proof that the updated version (including any new replacement mechanisms) is safe. To verify this system, we need a reflection principle to conclude that any updated version is safe assuming only that it was proved safe before the update was made. Furthermore, if we want to leave open the possibility that later updated versions retain a general replacement mechanism, we need a reflection principle that can be iterated.

The reflection principles in the previous section do not iterate. Using the assumption that a strongly inaccessible cardinal exists and is larger than the type of individuals they let us prove results of the form “If  $\ulcorner \varphi \urcorner$  is provable, then  $\varphi$ ” for every  $\varphi$  not containing type variables. However, the proof of  $\ulcorner \varphi \urcorner$  may only make use of the ordinary axioms of HOL: in particular, it cannot in turn assume that there is a strongly inaccessible cardinal. It is tempting to choose  $\varphi$  to be the formula `strongly_inaccessible  $\mathcal{U}(:\mu) \implies \psi$` , for some formula  $\psi$ . But this contains the type variable  $\mu$ , leading our automation to produce the unsatisfiable assumption  `$\mathcal{U}(:\mu) \prec \mathcal{U}(:\mu)$` . We could instead try to use a different type variable, as in `strongly_inaccessible  $\mathcal{U}(:\nu) \implies \psi$` ; this would lead to a theorem showing  $\psi$  under the assumptions that `strongly_inaccessible  $\mathcal{U}(:\mu)$` , `strongly_inaccessible  $\mathcal{U}(:\nu)$` ,  `$\mathcal{U}(:\nu) \prec \mathcal{U}(:\mu)$` , and the provability of  $\ulcorner \text{strongly\_inaccessible } \mathcal{U}(:\nu) \implies \psi \urcorner$ . However, while the proof in inner HOL may now assume that there is *one* inaccessible ( $\nu$ ), the theorem in outer HOL now assumes that there are *two* strongly inaccessible cardinals ( $\nu$  and  $\mu$ ), one of which is larger than the other.

Indeed, there must always be a stronger assumption in outer HOL than in inner HOL, since by Gödel’s second incompleteness theorem, no consistent proof system as strong as Peano Arithmetic can prove the reflection principle for itself: Choosing  $\varphi \equiv \text{F}$ , the identically false proposition, the assertion “If  $\ulcorner \text{F} \urcorner$  is provable, then  $\text{F}$ ” is equivalent to “ $\ulcorner \text{F} \urcorner$  is not provable”, which asserts consistency of the proof system. But by the second incompleteness theorem, a sufficiently strong proof system which can show its own consistency is inconsistent.

Nevertheless, the construction outlined above can be repeated any finite number of times, showing that if there are  $(n + 1)$  nested strongly inaccessible cardinals, then a proposition  $\varphi$  holds if it is provable under the assumption that there are  $n$  inaccessibles. Formally, we can define a term `LCA` of type `num  $\rightarrow$  ( $\mu \rightarrow \text{bool}$ )  $\rightarrow$  bool` such that `LCA 0  $\mathcal{U}(:\mu)$`  indicates that *ind* fits inside  $\mu$ , and such that `LCA (Suc  $t$ )  $\mathcal{U}(:\mu)$`  indicates that  $\mu$  is a strongly inaccessible cardinal and there is a strictly smaller subset  $Q$  of  $\mathcal{U}(:\mu)$  which satisfies `LCA  $t$   $Q$` :

$$\begin{aligned} \text{LCA } 0 \ P &\iff \mathcal{U}(:\text{ind}) \preceq P \\ \text{LCA (Suc } n) \ P &\iff \text{strongly\_inaccessible } P \wedge \exists Q. Q \subseteq P \wedge Q \prec P \wedge \text{LCA } n \ Q \end{aligned}$$

Then, we can show a reflection principle of the form “If  $\ulcorner \text{LCA } \bar{t} \ \mathcal{U}(:\mu) \implies \varphi \urcorner$  is provable, then `LCA (Suc  $t$ )  $\mathcal{U}(:\mu) \implies \varphi$` ”. This approach systematises the idea of allowing stronger systems to reflect the reasoning of weaker systems. As described below, we can strengthen this principle so that it requires proofs not about particular numerals  $\ulcorner \bar{t} \urcorner$  but about a universally quantified variable  $t$ .

To see the relevance to problems like verifying the self-replacing operating system, consider interpreting `Safe  $t$   $v$`  as “candidate system  $v$ , and any updated versions it permits, behave safely for  $t$  updates”, hence we want to show  $\forall t. \text{Safe } t \ v_0$  for the initial system  $v_0$ . Suppose we have proved `Safe 0  $v_0$` , that is, we have verified the initial system except for its replacement mechanism, and suppose the initial replacement mechanism requires a proof of the proposition  $\ulcorner \forall t. \text{LCA } t \ \mathcal{U}(:\mu) \implies \text{Safe } t \ \bar{v} \urcorner$  before installing update  $v$ . We need, for all  $t$ , to

show  $\text{Safe} (\text{Suc } t) v_0$ , but this is equivalent to showing  $\text{Safe } t v$  for all  $v$  that satisfy the replacement mechanism’s proof requirement. The following reflection principle enables us to do just that, provided we assume  $\text{LCA} (\text{Suc } t) \mathcal{U}(:\mu)$ .

$$\begin{aligned} & \vdash \forall v. \\ & (\text{thyof inner\_ctxt, []}) \vdash \lceil \forall t. \text{LCA } t \mathcal{U}(:\mu) \Rightarrow \text{Safe } t \bar{v} \rceil \Rightarrow \\ & \forall t. \text{LCA} (\text{Suc } t) \mathcal{U}(:\mu) \Rightarrow \text{Safe } t v \end{aligned}$$

Furthermore, the same reflection principle can be used within the proof required by the replacement mechanism, because the increase in the value of  $t$  passed to  $\text{LCA}$  is cancelled out by the decrease in the argument to  $\text{Safe}$  when considering a new version.

To prove the theorem above (for any predicate  $\text{Safe } t v$  not containing type variables), our automation shows that for every natural number  $t$ ,  $\text{LCA} (\text{Suc } t) \mathcal{U}(:\mu)$  implies the existence of a model of inner HOL in which  $\lceil \text{LCA } t \mathcal{U}(:\mu) \rceil$  is true (with the inner variable  $t$  being interpreted to equal its outer value). Thus, if the implication in the assumption of the theorem is provable, then  $\lceil \text{Safe } t \bar{v} \rceil$  is true in this model as well, and by our link between inner and outer HOL, it follows that  $\text{Safe } t v$ .<sup>14</sup>

At the top level, we are still left with the assumption of  $\text{LCA} (\text{Suc } t) \mathcal{U}(:\mu)$  for every  $t$ . A model of HOL in which this assumption is true can easily be constructed, for example, in ZFC extended with the much stronger assumption that there is a *Mahlo cardinal* [13, Chapter 8], a strongly inaccessible cardinal  $\kappa$  such that there are exactly  $\kappa$  strongly inaccessibles  $\preceq \kappa$ . The assumption that a Mahlo cardinal exists is not uncommon in set-theoretical work, and has been used in studies of dependent type theory to justify inductive-recursive definitions [2].

## 9 Related work

Reflection principles of the form “If  $\lceil \varphi \rceil$  is provable, then  $\varphi$ ”, as well as the uniform reflection principle “If  $\lceil \varphi(\bar{n}) \rceil$  is provable, then  $\varphi(n)$ ”, have been studied by Turing [22] and Feferman [4]. These authors consider sequences of theories obtained by starting with a theory  $T_0$  (Peano Arithmetic, say), and repeatedly constructing theories  $T_{\alpha+1}$  by adding to  $T_\alpha$  the reflection principle for all theories  $T_\beta$ ,  $\beta \leq \alpha$ . One can extend these sequences transfinitely by letting  $T_\lambda := \bigcup_{\alpha < \lambda} T_\alpha$  at limit ordinals  $\lambda$ . Turing and Feferman showed that, in a rather technical sense, sequences constructed in this way prove every true sentence of arithmetic (see Franzén [5] for an introduction, including an explanation why this statement is not as strong as it may appear). Set theorists have been interested in reflection principles, provable in ZF, which show that a sentence  $\varphi$  is true if it is provable from a fixed finite subset of ZF [13].

In the interactive theorem proving community, interest in reflection principles has mainly come from the perspective of *computational* reflection [10], which

<sup>14</sup> Our proof constructs a potentially slightly different model of inner HOL for each value of  $t$ ; this is, roughly, the origin of the term *model polymorphism*.

justifies the use of an efficient decision procedure by proving that if the decision procedure declares a sentence to be true, the sentence is in fact provable. For example, Allen et al. [1] extend Nuprl with a *reflection rule*, an inference rule to infer  $\varphi$  from “ $\ulcorner\varphi\urcorner$  is provable”. To avoid inconsistency, they stratify their rule in a manner similar to the transfinite progressions of Turing and Feferman: each invocation of the reflection rule is annotated with a level,  $\ell$ , and an invocation at level  $\ell$  requires a proof that  $\ulcorner\varphi\urcorner$  is provable using the reflection rule only at levels  $< \ell$ . Perhaps the most systematic use of reflection to justify more complex inference rules from simpler ones is in Milawa [18], which allows its entire proof checker to be replaced by a new version when given a proof that all sentences which are provable according to the new proof checker were also provable according to the old one.

Harrison [10], reviewing a large number of arguments and proposals for computational reflection, finds no evidence that it ever makes an otherwise infeasible proof technique feasible, though he concedes that in some cases, the speed-up can be significant. More recently, Coq’s `ssreflect` library [8] for computational reflection has been instrumental in the formal verification of the four-color theorem [7].

The algorithm in Section 4 bears a strong resemblance to the proof-producing translation (or code generation) algorithm presented by Myreen and Owens [19]. The input term in their algorithm is also an outer HOL term, but the target, instead of inner HOL, is the functional programming language CakeML. The semantics of CakeML is very different from the semantics of inner HOL, but the overall approach of producing certificate theorems bottom-up works similarly.

## 10 Conclusion

We have described automation for proving reflection principles of the form “If  $\ulcorner\varphi\urcorner$  is provable, then  $\varphi$ ” from assumptions about the existence of large cardinals. Based on this work, we have discussed an implementation of *model polymorphism* [3], which allows reflective reasoning to be used in the verification of self-replacing systems.

In this paper, we have focused on the automation for proving reflection principles. In future work, we plan to apply this to an implementation of a self-modifying, self-verifying system. It would also be interesting to apply these techniques to create an extensible version of HOL, similar to Milawa [18]; in particular, where Milawa requires a proof that a new version of the proof checker is *conservative* (i.e., only accepts proofs of propositions that were also provable according to the old proof checker), our semantic approach would allow us to instead require a proof that the new version is *sound* (i.e., whenever it accepts a proof of a proposition, that proposition is semantically true). It would be interesting to explore how these syntactic and semantic extension principles compare in practice.

*Acknowledgements* We thank Magnus Myreen for feedback on a draft of this paper. We also thank the anonymous reviewers for their helpful criticism.

## References

1. Allen, S.F., Constable, R.L., Howe, D.J., Aitken, W.E.: The semantics of reflected proof. In: Proc. LICS. pp. 95–105. IEEE Computer Society (1990)
2. Dybjer, P., Setzer, A.: A finite axiomatization of inductive-recursive definitions. In: Girard, J. (ed.) Proc. TCLA. LNCS, vol. 1581, pp. 129–146. Springer (1999)
3. Fallenstein, B., Soares, N.: Vingean reflection. Tech. rep., Machine Intelligence Research Institute, Berkeley, CA (2015)
4. Feferman, S.: Transfinite recursive progressions of axiomatic theories. *J. Symb. Log.* 27(3), 259–316 (1962)
5. Franzén, T.: Transfinite progressions: a second look at completeness. *B. Symb. Log.* 10(3), 367–389 (2004), <http://www.math.ucla.edu/~asl/bsl/1003/1003-003.ps>
6. Gödel, K.: Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik* 38(1), 173–198 (1931)
7. Gonthier, G.: The four colour theorem: Engineering of a formal proof. In: Kapur, D. (ed.) ASCM. LNCS, vol. 5081, p. 333. Springer (2007)
8. Gonthier, G., Mahboubi, A.: An introduction to small scale reflection in Coq. *J. Form. Reasoning* 3(2), 95–152 (2010)
9. Gordon, M.: From LCF to HOL: a short history. In: Plotkin, G.D., Stirling, C., Tofte, M. (eds.) *Proof, Language, and Interaction, Essays in Honour of Robin Milner*. pp. 169–186. The MIT Press (2000)
10. Harrison, J.: Metatheory and reflection in theorem proving: A survey and critique. Tech. Rep. CRC-053, SRI, Cambridge, UK (1995), <http://www.cl.cam.ac.uk/~jrh13/papers/reflect.dvi.gz>
11. Harrison, J.: Towards self-verification of HOL Light. In: Furbach, U., Shankar, N. (eds.) Proc. IJCAR. LNCS, vol. 4130, pp. 177–191. Springer (2006)
12. Harrison, J.: HOL Light: An overview. In: Berghofer, S., Nipkow, T., Urban, C., Wenzel, M. (eds.) Proc. TPHOLs. LNCS, vol. 5674, pp. 60–66. Springer (2009)
13. Jech, T.: Set theory. The third millennium edition, revised and expanded. Springer Monographs in Mathematics, Springer (2003)
14. Klein, G., Gamboa, R. (eds.): Proc. ITP, LNCS, vol. 8558. Springer (2014)
15. Kumar, R., Arthan, R., Myreen, M.O., Owens, S.: HOL with definitions: Semantics, soundness, and a verified implementation. In: Klein and Gamboa [14], pp. 308–324
16. Kumar, R., Arthan, R., Myreen, M.O., Owens, S.: Self-formalisation of higher-order logic. *J. Autom. Reasoning* (2015), submitted. Preprint at <https://cakeml.org>.
17. Mohamed, O.A., Muñoz, C.A., Tahar, S. (eds.): Proc. TPHOLs, LNCS, vol. 5170. Springer (2008)
18. Myreen, M.O., Davis, J.: The reflective Milawa theorem prover is sound - (down to the machine code that runs it). In: Klein and Gamboa [14], pp. 421–436
19. Myreen, M.O., Owens, S.: Proof-producing translation of higher-order logic into pure and stateful ML. *J. Funct. Program.* 24(2-3), 284–315 (2014)
20. Norrish, M., Huffman, B.: Ordinals in HOL: transfinite arithmetic up to (and beyond)  $\omega_1$ . In: Blazy, S., Paulin-Mohring, C., Pichardie, D. (eds.) Proc. ITP. LNCS, vol. 7998, pp. 133–146. Springer (2013)
21. Slind, K., Norrish, M.: A brief overview of HOL4. In: Mohamed et al. [17], pp. 28–32
22. Turing, A.M.: Systems of logic based on ordinals. *Proc. LMS* 2(1), 161–228 (1939)
23. Wenzel, M., Paulson, L.C., Nipkow, T.: The Isabelle framework. In: Mohamed et al. [17], pp. 33–38