

Safety and Conservativity of Definitions in HOL and Isabelle/HOL

ONDŘEJ KUNČAR, Technische Universität München

ANDREI POPESCU, Middlesex University London and Institute of Mathematics Simion Stoilow of the Romanian Academy

Definitions are traditionally considered to be a safe mechanism for introducing concepts on top of a logic known to be consistent. In contrast to arbitrary axioms, definitions should in principle be treatable as a form of abbreviation, and thus compiled away from the theory without losing provability. In particular, definitions should form a conservative extension of the pure logic. These properties are crucial for modern interactive theorem provers, since they ensure the consistency of the logic, as well as a valid environment for total/certified functional programming.

We prove these properties, namely, safety and conservativity, for Higher-Order Logic (HOL), a logic implemented in several mainstream theorem provers and relied upon by thousands of users. Some unique features of HOL, such as the requirement to give non-emptiness proofs when defining new types and the impossibility to unfold type definitions, make the proof of these properties, and also the very formulation of safety, nontrivial.

Our study also factors in the essential variation of HOL definitions featured by Isabelle/HOL, a popular member of the HOL-based provers family. The current work improves on recent results which showed a weaker property, consistency of Isabelle/HOL's definitions.

CCS Concepts: • **Theory of computation** → **Logic and verification; Higher order logic; Type structures; Interactive proof systems;**

Additional Key Words and Phrases: higher-order logic (HOL), interactive theorem proving, type definitions, conservative extensions, Isabelle/HOL

ACM Reference format:

Ondřej Kunčar and Andrei Popescu. 2017. Safety and Conservativity of Definitions in HOL and Isabelle/HOL. *Proc. ACM Program. Lang.* 1, 1, Article 1 (January 2017), 28 pages.
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Higher-Order Logic (HOL) [Pitts 1993] (recalled in Section 3 of this paper) is an important logic in the theorem proving community. It forms the basis of several interactive theorem provers (also known as proof assistants), including HOL4 [Gordon and Melham 1993; Slind and Norrish 2008], HOL Light [Harrison 1996], Isabelle/HOL [Nipkow et al. 2002], ProofPower-HOL [Arthan 2004] and HOL Zero [Adams 2010].

In addition to supporting the development of formalized mathematics, most modern interactive theorems provers also include a functional programming language, supporting the paradigm of *total programming* [Turner 2004]. For example, in provers based on type theory such as Agda [Bove, Dybjer, and Norell Bove et al.], Coq [Bertot and Casteran 2004] and Matita [Asperti et al. 2011], totality is ensured by a global *strong normalization property*. There is a tight relationship between this property, allowing functions/programs to be reduced to a normal form by recursively unfolding all definitions and reducing all redexes, and the logical consistency of these systems.

2017. 2475-1421/2017/1-ART1 \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

In HOL-based provers, programming is supported by a different mechanism: All recursive datatype specifications and all recursive specifications of functions on these datatypes are translated into *nonrecursive* HOL primitives, i.e., constant and type definitions; then the recursive specifications are proved automatically as *theorems* in the logic. This scheme involves a massive background compilation and proof process (supported by tools consisting of tens of thousands of lines of code, e.g., [Blanchette et al. 2014; Krauss 2009; Melham 1989]). It ensures a high degree of trustworthiness—given that all constructions must pass through the “firewall” of HOL’s minimalistic kernel. In particular, a potential bug in the compilation tools could cause correct user specifications to fail, but will not introduce logical inconsistencies unless the kernel has a bug.

In this paper, we turn our attention to the HOL kernel itself, which is the warrant of logical consistency and certified programming in the above scheme. In spite of extensive foundational studies and the relative simplicity of the logic, the normalization process underlying the HOL kernel, i.e., the process of unfolding the HOL definitions, remains less understood than its corresponding process in type theory, and occasionally leads to controversial design decisions and heated debates—as we are about to show, after recalling some background information.

While its ideas go back a long way (to the work of Alonzo Church [Church 1940] and beyond), HOL contains a unique blend of features proposed by Mike Gordon at the end of the eighties, inspired by practical verification needs: Its type system is the rank-one polymorphic extension of simple types, generated using the function-space constructor from two base types, `bool` and `ind`; its terms have built-in equality (from which all the usual connectives and quantifiers are derived); deduction, operating on terms of type `bool` called formulas, is regulated by the built-in axioms of Equality, (Hilbert) Choice and Infinity (for the type `ind`). In addition to this purely logical layer, which we shall refer to as *initial HOL*, users can perform constant and type declarations and definitions. Type definitions proceed by indicating a predicate on an existing type and carving out the new type from the subset satisfying the predicate. For accepting a type definition, the system requires a proof that the subset is nonempty (the predicate has a witness). This is because *HOL types are required to be nonempty*—a major design decision, with practical and theoretical ramifications [Gordon and Melham 1993; Paulson 1990]. No new axioms are accepted (more precisely, they are strongly discouraged), besides the aforementioned definitions. This minimalistic, *definitional approach* offers good protection against the accidental introduction of *inconsistency* (the possibility to prove False).

Isabelle/HOL [Nipkow et al. 2002] is a notable member of the HOL family, and a maverick to some extent. It implements an essential variation of HOL, where constant definitions can be overloaded in an ad hoc manner, for different instances of their types. This flexibility forms the basis of Haskell-style type classes [Nipkow and Snelting 1991],¹ a feature that allows for lighter, suppler formalizations and should probably be credited, together with the high-level structured proof language [Wenzel 1999], the powerful automation [Paulson 2010] and the convenient user interface [Wenzel 2014], for Isabelle/HOL’s wide popularity and prolificness: thousands of users in both academia and industry, a large library of formalized results [Isabelle 2016; Klein et al. 2016], major verification success stories [Esparza et al. 2013; Klein et al. 2010; Lochbihler 2010; Paulson 2015].

The founding fathers of HOL have paid special attention to consistency and related properties. Andrew Pitts designed a custom notion of *standard model* [Pitts 1993], aimed at smoothly accommodating both polymorphism and type definitions. He proved that constant and type definitions are *model-theoretically conservative w.r.t. standard models*: Any standard model of a theory can be

¹Type classes do not require any additional extension of the logic, but are completely reduced (including at the level of proofs) to HOL with type definitions and ad hoc overloaded constants [Wenzel 1997, Section 5].

expanded to a standard model of the theory plus the definitions. This of course implies consistency of HOL with definitions. Surprisingly, the founding fathers have not looked into the more customary notion of *proof-theoretic conservativity*, which we shall simply call *conservativity*. It states that, by adding new constants and types and their definitions, nothing new can be proved in the old language. This does not follow from the model-theoretic version (because of the restriction to *standard* models, for which deduction is not complete). In fact, as we discuss below, it does not even hold in general.

In Isabelle/HOL, the foundational problem is more challenging. Here, even the *consistency* of definitions has not been fully understood until very recently (Section 2.2). The culprit is precisely the feature that contributes to Isabelle/HOL’s popularity—ad hoc overloading—which has a delicate interaction with type definitions [Kunčar and Popescu 2015, Section 1].

Motivated by the desire to settle the Isabelle foundations, early work by Wenzel formulates criteria for safety of definitions in HOL-like logics [Wenzel 1997]. For a theory extension $\Theta_1 \subseteq \Theta_2$, he considers (proof-theoretic) conservativity, a property much stronger than preservation of consistency, to be a minimum requirement for deeming a theory extension truly definitional [Wenzel 1997, p.7]. In fact, he argues for an even stronger notion, *meta-safety*. Let Σ_1 and Σ_2 be the languages (signatures) of Θ_1 and Θ_2 , respectively. (Thus, $\Sigma_1 \subseteq \Sigma_2$.) Meta-safety requires that, whenever a Σ_2 -formula φ is deducible from Θ_2 , there exists a Σ_1 -formula $\varphi[\dots, t/c, \dots]$, obtained by replacing all the items $c \in \Sigma_2 \setminus \Sigma_1$ with some suitable Σ_1 -terms t , which is deducible from Θ_1 . This way, the items c can be considered to be “defined” because they can always be compiled away without losing provability. He also shows that, under appropriate well-formedness restrictions, a set of constant definitions, even overloaded as in Isabelle/HOL, forms a meta-safe extension.

However, as formulated, meta-safety does not apply to type definitions, because in HOL it is impossible to replace a defined type with its defining expression. In fact, Wenzel makes the following observation: *In general, type definitions in HOL are not even consistency-preserving, let alone conservative (let alone meta-safe in any reasonable way)*, as witnessed by the following example. Consider the HOL theory consisting of a single formula φ stating that no type has *precisely* three elements (i.e. for all types α , if α has at most three elements x, y, z then two of them must be equal):

$$\forall x, y, z : \alpha. (\forall v : \alpha. v = x \vee v = y \vee v = z) \longrightarrow x = y \vee x = z \vee y = z$$

The theory $\{\varphi\}$ is consistent since there exists a model that satisfies it—the full-frame model of initial HOL, where all finite types are function-space combinations over `bool`, hence their cardinality is a power of two, in particular, no type has cardinality three. On the other hand, the extension of $\{\varphi\}$ with the definition of a type having three elements, $\tau = \{0, \text{Suc } 0, \text{Suc}(\text{Suc } 0)\}$, is clearly inconsistent—which exhibits a type definition that does not preserve consistency. This analysis has led Wenzel, who is Isabelle’s long-standing lead developer and release manager, to deem type definitions *axiomatic* (i.e., having *zero* consistency or conservativity guarantees attached) rather than definitional. This departure from a well-established HOL tradition has generated confusion and misunderstanding amongst Isabelle/HOL’s users and developers [Wolff 2015].

But the above counterexample involves a non-definitional theory— φ is not a definition, but merely an axiom that happens to be consistent. Thus, the counterexample only shows that, unlike constant definitions, type definitions do not preserve consistency, a fortiori, are not conservative, *over an arbitrary (axiomatic) theory*. Nonetheless, it is still legitimate to ask:

Are arbitrary combinations of constant and type definitions conservative over initial HOL?

And are they even meta-safe (again, over initial HOL) in a suitable sense?

We believe these are important questions for deepening our understanding of the nature of HOL and Isabelle/HOL definitions. Conservativity also provides the most compelling way of witnessing

	Over Initial HOL	Over Arbitrary HOL Theories	
Constant Definitions	Yes (from the right)	Yes [Wenzel 1997]	Yes (from below)
Constant Definitions Mixed with Type Definitions	Yes (this paper)	No [Wenzel 1997]	Yes [Pitts 1993]
Isabelle-HOL Constant Definitions	Yes (from the right)	Yes [Wenzel 1997]	No [Wenzel 1997]
Isabelle-HOL Constant Definitions Mixed with Type Definitions	Yes (this paper)	No [Wenzel 1997]	No (from above)
	(Proof-Theoretic) Conservativity		Model-Theoretic Conservativity w.r.t. Standard Models

Fig. 1. Conservativity of Definitions in HOL and Isabelle/HOL

consistency: Any proof of False using definitions can be traced down to a proof of False in initial HOL (the latter being manifestly consistent thanks to its standard set-theoretic semantics). This is especially relevant for the brittle foundational terrain of Isabelle/HOL, where it should help rehabilitating type definitions as genuine, safe definitions.

In this paper, we provide a positive answer to both questions. Figure 1 shows our conservativity results in the context of similar known facts. (Note that, for Isabelle/HOL constant definitions, proof-theoretic conservativity holds, roughly because all definitions can be unfolded; by contrast, model-theoretic conservativity fails, due to ad hoc overloading—a declared constant is allowed to have any interpretation in a model, but this can be in conflict with a future definition of an instance of that constant.)

First, we focus on traditional HOL, where we formulate meta-safety by defining translation operators for types and terms that unfold the definitions (Section 4). Unfolding a type definition has to be done in an indirect fashion, since HOL does not support comprehension/refinement types (of the form $\{x : \sigma \mid t \ x\}$). Namely, a formula operating on defined types will be relativized to a formula on the original, built-in types that hosted the type definitions; so the “unfolding” of a defined type will be a predicate on its host type. Since type definitions are paired with nonemptiness proofs (in the current contexts, having available all the previously introduced definitions), we are forced to proceed gradually, one definition at a time. Consequently, the proof of meta-safety (also leading to conservativity) is itself gradual, in a feedback loop between preservation of deduction, commutation with substitution, and nonemptiness of the relativization predicates.

We organized the proof development for traditional HOL modularly, separating lemmas about termination of the definitional dependency relation. This allows a smooth upgrade to the more complex case of Isabelle/HOL (Section 5), where termination is no longer ensured by the historic order of definitions, but via a more global approach. Due to ad hoc overloading, here the translations no longer commute with type substitution. We recover from this “anomaly” by mining the proofs and weakening the commutation lemma—leading to an Isabelle/HOL version of the results.

In the appendix, we give more details on the HOL logic concepts and show some omitted proofs. Moreover, we implemented for Isabelle/HOL the unfolding and relativization functions presented in this paper, and used them to check the paper’s examples. The documented implementation is provided as additional material.

2 MORE RELATED WORK

There is a vast literature on the logical foundations of theorem provers, which we will not attempt to survey here. We focus on work that is directly relevant to our present contribution, from the point of view of either the object logic or the techniques used.

2.1 HOL Foundations

Wiedijk [Wiedijk 2009] defines *stateless HOL*, a version of HOL where terms and types carry *in their syntax* information about the defined constants and type constructors. Kumar et al. [Kumar et al. 2014] define a set-theoretic (Pitts-style) model for stateless HOL and a translation from standard (stateful) HOL with definitions to stateless HOL, thus proving the consistency of both. Their stateful to stateless HOL translation is similar to our translation, in that they both internalize the definitions (which are part of “the state”) into “stateless” formulas; however, for conservativity, we need to appeal to pure HOL entities, not to syntactically enriched ones. In a subsequent paper [Kumar et al. 2016], the same authors renounce the stateless HOL detour and prove model-theoretic conservativity directly on initial HOL.

Kumar et al.’s work, which has been mechanized in the HOL4 theorem prover, is based on pioneering self-verification work by Harrison [Harrison 2006], who uses HOL Light to give semantic proofs of soundness of the HOL logic without definitional mechanisms, in two flavors: either after removing the infinity axiom from the object HOL logic, or after adding a “universe” axiom to HOL Light.

2.2 Isabelle/HOL Foundations

Wenzel’s work cited in the introduction [Wenzel 1997] proved meta-safety and conservativeness of constant definitions but left type definitions aside. In spite of Wenzel’s theoretical observation that orthogonality and termination are required to ensure meta-safety, overloading of constants remained unchecked in Isabelle/HOL for many years—until Obua [Obua 2006] looked into the problem and proposed a way to implement Wenzel’s observation with an external termination checker. Obua also aimed to extend the scope of consistency by factoring in type definitions. But his syntactic proof missed out possible inconsistencies through delayed overloading intertwined with type definitions. Soon after, Wenzel designed and implemented a more structural solution based on work of Haftmann, Obua and Urban (parts of which are reported in [Haftmann and Wenzel 2006]).

The foundational work on Isabelle/HOL was resumed in 2014, after the aforementioned inconsistencies caused by delayed overloading and type definitions were discovered. To address the problem, Kunčar and Popescu [Kunčar and Popescu 2015] defined a new dependency relation, operating on constants *and types* (which became part of the system starting from Isabelle2016). Employing a nonstandard semantics, they proved that, after these modifications, any definitional theory is consistent. In more recent work, the same authors gave an alternative syntactic proof, based on translating HOL to a richer logic, HOLC, having comprehension types as first-class citizens [Kunčar and Popescu 2017]. The current paper improves on these results, by proving properties much stronger than consistency.

2.3 Other Work

The interactive theorem proving community can be roughly divided in two successful camps: provers based on type theory (Agda, Coq, Matita, etc.) and provers based on HOL. For the former, the notion of normalizing terms is fairly well studied and well understood [Abel et al. 2007; Altenkirch 1993; Barras 2010; Coquand et al. 1990; Coquand and Spiwack 2006; Geuvers 1993]. Our notion of meta-safety can be seen as the HOL counterpart of type-theoretic normalization, hence as a foundation for HOL-based programming. Of course, the technical challenges we face in HOL are quite different—here, it is not the expressiveness of the logic or of its underlying type system (e.g., fancy dependent types or polymorphism) that complicates the argument, but to a large extent its *lack of expressiveness*: The logic disallows unfolding type definitions, which forces us into a labyrinth of relativization techniques. Another difference is that HOL is an inherently classical

logic: Type definitions require possibly non-constructive proofs of nonemptiness, and the Hilbert Choice is paramount. This makes our proof translations associated to normalization less clean than in type theory.

Other foundational work for theorem provers includes Myreen and Davis’s mechanized proof of consistency for Milawa [Myreen and Davis 2014], a prover based on first-order logic in the style of ACL2 [Kaufmann et al. 2000], and Owre and Shankar’s set-theoretic semantics of PVS [Owre and Shankar 1999]—featuring a logic similar to HOL, but different by the presence of dependent types and the absence of polymorphism.

Outside the world of theorem proving, conservative extensions are widely employed in mathematical logic, e.g., in the very popular Henkin technique for proving completeness [Henkin 1949]. They are also employed in algebraic specifications to achieve desirable modularity properties [Sannella and Tarlecki 2012]. However, in these fields, *definitional* extensions are often trivially conservative, thanks to their simple equational structure and freshness conditions.

3 HOL PRELIMINARIES

By HOL, we mean classical higher-order logic with Infinity, Choice and rank-one polymorphism, and mechanisms for constant and type definitions and declarations. This section explains all these concepts and features in detail.

3.1 Syntax

All throughout this paper, we fix the following:

- an infinite set TVar , of *type variables*, ranged by α, β
- an infinite set VarN , of (*term*) *variable names*, ranged by x, y, z

A *type structure* is a pair (K, arOf) where:

- K is a set of symbols, ranged by k , called *type constructors*, containing three special symbols: “bool”, “ind” and “ \Rightarrow ” (aimed at representing the type of booleans, an infinite type of individuals and the function type constructor, respectively)
- $\text{arOf} : K \Rightarrow \mathbb{N}$ is a function associating arities to the type constructors, such that $\text{arOf}(\text{bool}) = \text{arOf}(\text{ind}) = 0$ and $\text{arOf}(\Rightarrow) = 2$.

The *types* associated to (K, arOf) , ranged by σ, τ , are defined as follows:

$$\sigma ::= \alpha \mid (\sigma_1, \dots, \sigma_{\text{arOf}(k)}) k$$

Thus, a type is either a type variable or an n -ary type constructor k postfix-applied to a number of types corresponding to its arity. We write $\text{Type}_{(K, \text{arOf})}$ for the set of types associated to (K, arOf) .

A *signature* is a tuple $\Sigma = (K, \text{arOf}, \text{Const}, \text{tpOf})$, where:

- (K, arOf) is a type structure
- Const , ranged over by c , is a set of symbols called *constants*, containing four special symbols: “=”, “ ε ”, “zero” and “suc” (aimed at representing equality, Hilbert choice of some element from a type, zero and successor, respectively)
- $\text{tpOf} : \text{Const} \Rightarrow \text{Type}$ is a function associating a type to every constant, such that:

$$\begin{aligned} \text{tpOf}(\longrightarrow) &= \text{bool} \Rightarrow \text{bool} \Rightarrow \text{bool} & \text{tpOf}(=) &= \alpha \Rightarrow \alpha \Rightarrow \text{bool} \\ \text{tpOf}(\varepsilon) &= (\alpha \Rightarrow \text{bool}) \Rightarrow \alpha & \text{tpOf}(\text{zero}) &= \text{ind} & \text{tpOf}(\text{suc}) &= \text{ind} \Rightarrow \text{ind} \end{aligned}$$

For the rest of this section, we fix a signature $\Sigma = (K, \text{arOf}, \text{Const}, \text{tpOf})$. We usually write Type_Σ , or simply Type , instead of $\text{Type}_{(K, \text{arOf})}$.

$\text{TV}(\sigma)$ is the set of type variables of a type σ . A *type substitution* is a function $\rho : \text{TVar} \Rightarrow \text{Type}$. We let TSubst denote the set of type substitutions. The application of ρ to a type σ , written $\sigma[\rho]$, is

defined recursively by $\alpha[\rho] = \rho(\alpha)$ and $((\sigma_1, \dots, \sigma_m) k)[\rho] = (\sigma_1[\rho], \dots, \sigma_m[\rho]) k$. If $\alpha_1, \dots, \alpha_m$ are all different, we write $\tau_1/\alpha_1, \dots, \tau_n/\alpha_m$ for the type substitution that sends α_i to τ_i and each $\beta \notin \{\alpha_1, \dots, \alpha_m\}$ to β . Thus, $\sigma[\tau_1/\alpha_1, \dots, \tau_n/\alpha_m]$ is obtained from σ by substituting, for each i , τ_i for all occurrences of α_i .

We say that σ is an *instance* of τ via ρ , written $\sigma \leq_\rho \tau$, if $\tau[\rho] = \sigma$. We say that σ is an *instance* of τ , written $\sigma \leq \tau$, if there exists $\rho \in \text{TSubst}$ such that $\sigma \leq_\rho \tau$. Two types σ_1 and σ_2 are called *orthogonal*, written $\sigma_1 \# \sigma_2$, if they have no common instance; i.e., for all τ it holds that $\tau \not\leq \sigma_1$ or $\tau \not\leq \sigma_2$.

Given $\rho_1, \rho_2 \in \text{TSubst}$, we write $\rho_1 \cdot \rho_2$ for their *composition*, defined as $(\rho_1 \cdot \rho_2)(\alpha) = (\rho_1(\alpha))[\rho_2]$. It is easy to see that, for all types σ , it holds that $\sigma[\rho_1 \cdot \rho_2] = \sigma[\rho_1][\rho_2]$.

A (*typed*) *variable* is a pair of a variable name x and a type σ , written x_σ . We let Var denote the set of variables. A *constant instance* is a pair of a constant and a type, written c_σ , such that $\sigma \leq \text{tpOf}(c)$. We let CInst denote the set of constant instances. We extend the notions of being an instance (\leq) and being orthogonal ($\#$) from types to constant instances:

$$c_\tau \leq d_\sigma \text{ iff } c = d \text{ and } \tau \leq \sigma \qquad c_\tau \# d_\sigma \text{ iff } c \neq d \text{ or } \tau \# \sigma$$

The signature's *terms*, ranged over by s, t , are defined by the grammar:

$$t ::= x_\sigma \mid c_\sigma \mid t_1 t_2 \mid \lambda x_\sigma. t$$

Thus, a term is either a variable, or a constant instance, or an application, or an abstraction. As usual, we identify terms modulo alpha-equivalence. We let Term_Σ , or simply Term , ranged by s and t , denote the set of terms. Typing is defined as a binary relation between terms and types, written $t : \sigma$, inductively as follows:

$$\frac{x_\sigma \in \text{Var}}{x_\sigma : \sigma} \qquad \frac{c_\sigma \in \text{CInst}}{c_\sigma : \sigma} \qquad \frac{t_1 : \sigma \Rightarrow \tau \quad t_2 : \sigma}{t_1 t_2 : \tau} \qquad \frac{t : \tau}{\lambda x_\sigma. t : \sigma \Rightarrow \tau}$$

We can apply a type substitution ρ to a term t , written $t[\rho]$, by applying it to the types of all variables and constant instances occurring in t with the usual renaming of bound variables if they get captured. $\text{FV}(t)$ is the set of t 's free variables. The term t is called *closed* if it has no free variables: $\text{FV}(t) = \emptyset$. We write $t[s/x_\sigma]$ for the term obtained from t by capture-free substituting the term s for all free occurrences of x_σ .

A *formula* is a term of type bool . The formula connectives (e.g., \wedge and \longrightarrow) and quantifiers (\forall and \exists) are defined in the usual way, starting from the equality primitive—the appendix gives details. The if-then-else construct, if_t_e , is defined as follows, given $b : \text{bool}$, $t_1 : \sigma$ and $t_2 : \sigma$

$$\text{if_t_e } b \ t_1 \ t_2 = \varepsilon (\lambda x_\sigma. (b \longrightarrow x_\sigma = t_1) \wedge (\neg b \longrightarrow x_\sigma = t_2))$$

Its behavior is the expected one: it equals t_1 if b is True and equals t_2 if b is False .

To avoid confusion with the object-logic definitions that we discuss later, we will treat the logical connectives and quantifiers and the if-then-else operator as mere abbreviations (i.e., meta-level definitions of certain HOL terms). When writing terms, we sometimes omit the types of variables if they can be inferred. For example, we write $\lambda x_\sigma. x$ instead of $\lambda x_\sigma. x_\sigma$. A *theory* (over Σ) is a set of closed (Σ -)formulas.

3.2 Axioms and Deduction

The HOL axioms, forming the set Ax , are the usual Equality axioms, the Infinity axioms (stating that suc is different from 0 and is injective, which makes the type ind infinite), the classical Excluded Middle and the Choice axiom, which states that the Hilbert choice operator returns an element satisfying its argument predicate (if nonempty): $p_{\alpha \Rightarrow \text{bool}} x \longrightarrow p (\varepsilon p)$.

A *context* Γ is a finite set of formulas. We write $\alpha \notin \Gamma$ to indicate that the type variable α does not appear in any formula in Γ ; similarly, $x_\sigma \notin \Gamma$ will indicate that x_σ does not appear *free* in any formula in Γ . We define *deduction* as a ternary relation \vdash between theories D , contexts Γ and formulas φ , written $D; \Gamma \vdash \varphi$.

$$\begin{array}{c}
\frac{}{D; \Gamma \vdash \varphi [\varphi \in \text{Ax} \cup D]} \text{(FACT)} \qquad \frac{}{D; \Gamma \vdash \varphi [\varphi \in \Gamma]} \text{(ASSUM)} \qquad \frac{D; \Gamma \vdash \varphi}{D; \Gamma \vdash \varphi[\sigma/\alpha]} \text{(T-INST)} \\
\\
\frac{D; \Gamma \vdash \varphi}{D; \Gamma \vdash \varphi[t/x_\sigma]} \text{(INST)} \qquad \frac{}{D; \Gamma \vdash (\lambda x_\sigma. t) s = t[s/x_\sigma]} \text{(BETA)} \qquad \frac{D; \Gamma \vdash f x_\sigma = g x_\sigma}{D; \Gamma \vdash f = g} \text{(EXT)} \quad [x_\sigma \notin \Gamma] \\
\\
\frac{D; \Gamma \cup \{\varphi\} \vdash \chi}{D; \Gamma \vdash \varphi \longrightarrow \chi} \text{(IMPL)} \qquad \frac{D; \Gamma \vdash \varphi \longrightarrow \chi \quad D; \Gamma \vdash \varphi}{D; \Gamma \vdash \chi} \text{(MP)}
\end{array}$$

The axioms and the deduction rules we gave here are (a variant of) the standard ones for HOL (as in, e.g., [Gordon and Melham 1993; Harrison 2006]). We write $D \vdash \varphi$ instead of $D; \emptyset \vdash \varphi$ and $\vdash \varphi$ instead of $\emptyset; \emptyset \vdash \varphi$ (that is, we omit empty contexts and theories). Note that the HOL axioms are not part of the parameter theory D , but are wired together with D in the (FACT) axiom. So $\vdash \varphi$ indicates that φ is provable from the HOL axioms only.

3.3 HOL Definitions and Declarations

Besides deduction, another main component of the HOL logic is a mechanism for introducing new constants and types by spelling out their definitions.

The *built-in type constructors* are `bool`, `ind` and \Rightarrow . The *built-in constants* are `=`, `ε` , `zero` and `suc`. Since the built-in items have an already specified behavior (by the HOL axioms), only non-built-in items can be defined.

DEFINITION 1.

Constant Definitions: Given a non-built-in constant c such that $\text{tpOf}(c) = \sigma$ and a closed term $t : \sigma$, we let $c_\sigma \equiv t$ denote the formula $c_\sigma = t$. We call $c_\sigma \equiv t$ a *constant definition* provided $\text{TV}(t) \subseteq \text{TV}(c_\sigma)$ (i.e., $\text{TV}(t) \subseteq \text{TV}(\sigma)$).

Type Definitions: Given types τ and σ and a closed term $t : \sigma \Rightarrow \text{bool}$, we let $\tau \equiv t$ denote the formula

$$\exists \text{rep}_{\tau \Rightarrow \sigma}. \text{One_One}_{\text{rep}} \wedge (\forall y_\sigma. t y \leftrightarrow (\exists x_\tau. y = \text{rep } x))$$

where $\text{One_One}_{\text{rep}}$ is the formula stating that rep is one-to-one (injective), namely, $\forall x_\tau, y_\tau. \text{rep } x = \text{rep } y \longrightarrow x = y$. We call $\tau \equiv t$ a *type definition*, provided τ has the form $(\alpha_1, \dots, \alpha_m) k$ such that k is a non-built-in type constructor, the α_i 's are all distinct type variables and $\text{TV}(t) \subseteq \{\alpha_1, \dots, \alpha_m\}$. (Hence, we have $\text{TV}(t) \subseteq \text{TV}(\tau)$, which also implies $\text{TV}(\sigma) \subseteq \text{TV}(\tau)$.)

A type definition expresses the following: The new type $(\alpha_1, \dots, \alpha_m) k$ is embedded in its host type σ via some one-to-one function rep , and the image of this embedding consists of the elements of σ for which t holds. Since types in HOL are required to be nonempty, the definition is only accepted if the user provides a proof that $\exists x_\sigma. t x$ holds. Thus, *to perform a type definition, one must give a nonemptiness proof*.

Type and Constant Declarations: Declarations in HOL are a logical extension mechanism which is significantly milder than definitions—they simply add new items to the signature as “uninterpreted,” without providing any definition.

3.4 Signature Extensions and the Initial Signature

In the remainder of this paper, when necessary for disambiguation, we will indicate the signature Σ as a subscript when denoting various sets and relations associated to it: Type_Σ , Term_Σ , CInst_Σ , \vdash_Σ , etc.

Given a signature $\Sigma = (\text{K}, \text{arOf}, \text{Const}, \text{tpOf})$ and an item u , we write $u \in \Sigma$ to mean that $u \in \text{K}$ or $u \in \text{Const}$. Given signatures $\Sigma = (\text{K}, \text{arOf}, \text{Const}, \text{tpOf})$ and $\Sigma' = (\text{K}', \text{arOf}', \text{Const}', \text{tpOf}')$, we say Σ is *included in* Σ' , or Σ' *extends* Σ , written $\Sigma \subseteq \Sigma'$, if $\text{K} \subseteq \text{K}'$, $\text{Const} \subseteq \text{Const}'$ and the functions arOf' and tpOf' are extensions of arOf and tpOf , respectively. We write $u \in \Sigma' \setminus \Sigma$ to mean $u \in \Sigma'$ and $u \notin \Sigma$. If $c \notin \text{Const}$ and $\sigma \in \text{Type}_\Sigma$, we write $\Sigma \cup \{(c, \sigma)\}$ for the extension of Σ with a new constant c of type σ . Similarly, if $k \notin \text{K}$, we write $\Sigma \cup \{(k, n)\}$ for the extension of Σ with a new type constructor k of arity n .

We write Σ_{init} for the *initial signature*, containing only built-in type constructors and constants. Note that, by definition, any signature extends the initial signature.

4 CONSERVATIVITY OF HOL DEFINITIONS

A HOL development, i.e., a session of interaction with the HOL logic from a user's perspective, consists of intertwining definitions, declarations and (statements and proofs of) theorems. Since theorems are merely consequences of definitions, we will not model them explicitly, but focus on definitions and declarations.

Let $\Sigma = (\text{K}, \text{arOf}, \text{Const}, \text{tpOf})$ be a signature and let D be a finite theory over Σ .

DEFINITION 2. D is said to be a *well-formed definitional theory* if $D = \{def_1, \dots, def_n\}$, where each def_i is a (type or constant) definition of the form $u_i \equiv t_i$, and there exist the signatures $\Sigma^1, \dots, \Sigma^n$ and $\Sigma_0, \Sigma_1, \dots, \Sigma_n$ such that $\Sigma_0 = \Sigma_{\text{init}}$, $\Sigma_n = \Sigma$ and the following hold for all $i \in \{1, \dots, n\}$:

- (1) $t_i \in \text{Term}_{\Sigma_i}$ and Σ_i is the extension of Σ^{i-1} with a fresh item defined by def_i , namely:
 - (1.1) If u_i has the form $(\alpha_1, \dots, \alpha_m) k$, then $k \notin \Sigma^{i-1}$ and $\Sigma_i = \Sigma^{i-1} \cup \{(k, m)\}$
 - (1.2) If u_i has the form c_σ , then $c \notin \Sigma^{i-1}$ and $\Sigma_i = \Sigma^{i-1} \cup \{(c, \sigma)\}$
- (2) If def_i is a type definition, meaning u_i is a type and $t_i : \sigma \Rightarrow \text{bool}$, then $\{def_1, \dots, def_{i-1}\} \vdash_{\Sigma^{i-1}} \exists x_\sigma. t_i x$
- (3) $\Sigma_{i-1} \subseteq \Sigma^i$

These conditions express that the theory D consists of intertwined definitions and declarations. The chain of extensions

$$\Sigma_{\text{init}} = \Sigma_0 \subseteq \Sigma^1 \subseteq \Sigma_1 \subseteq \Sigma^2 \subseteq \Sigma_2 \dots \subseteq \Sigma^n \subseteq \Sigma_n = \Sigma,$$

starting from the initial signature and ending with Σ , alternates sets of declarations (the items in $\Sigma^i \setminus \Sigma_{i-1}$) with definitions (the unique item u_i in $\Sigma_i \setminus \Sigma^{i-1}$ being defined by def_i , i.e., as $u_i \equiv t_i$). As shown by condition (2), in the case of type definitions, we also require proofs of non-emptiness of the defining predicate t (from the definitions available so far).

In short, the above conditions state something very basic: Definitions are introduced one at a time and the defined symbols are fresh. This is clearly obeyed by correct implementations of standard HOL, such as HOL4 and HOL Light. (By contrast, the Isabelle/HOL-specific conditions in Section 5 will involve the more complex notions of orthogonality and termination.)

DEFINITION 3. A theory E over Σ is said to be a (*proof-theoretic*) *conservative extension of initial HOL* if any formula proved from E that belongs to the initial signature Σ_{init} could have been proved without E or the types and constants from outside of Σ . Formally: For all $\varphi \in \text{Fmla}_{\Sigma_{\text{init}}}$, $E \vdash_\Sigma \varphi$ implies $\vdash_{\Sigma_{\text{init}}} \varphi$.

4.1 Roadmap

In what follows, we fix a well-formed definitional theory D and use for it the notations introduced in Def. 2, e.g., Σ , Σ_i . We first sketch the main ideas of our development, motivating the choice of the concepts. The more formal definitions and proofs will be given in the following subsections.

Our two main goals are to *formulate and prove D 's meta-safety* and to *prove D 's conservativity*. As with any respectable notion of its kind, meta-safety will easily yield conservativity, so we concentrate our efforts on the former.

Recall that, for a Σ -formula φ provable from D , meta-safety should allow us to replace all the defined items in φ with items in the initial signature without losing provability, i.e., obtaining a deducible Σ_{init} -formula φ' . For constants, the procedure is clear: Any defined constant c appearing in φ is replaced with its defining term t , then any defined constant d appearing in t is replaced with its defining term, and so on, until (hopefully) the process terminates and we are left with built-in items only.

But how about for types τ occurring in φ ? A HOL type definition $\tau \equiv t$ where $t : \sigma \Rightarrow \text{bool}$, is not an equality (there is no type equality in HOL), but a formula asserting the existence of a bijection between τ and the set of elements of Σ for which the predicate t holds. So it cannot be “unfolded.” First, let us make the simplifying assumption that $\sigma \in \text{Type}_{\Sigma_{\text{init}}}$ and $t \in \text{Term}_{\Sigma_{\text{init}}}$. Then the only reasonable Σ_{init} -substitute for τ is its host type σ ; however, after the replacement of τ by σ , the formula needs to be adjusted not to refer to the whole σ , but only to the isomorphic copy of τ —in other words, the formula needs to be relativized to the predicate t . In general, σ or t may themselves contain defined types or constants, which will need to be processed similarly, and so on, recursively. In summary:

- for each type τ , we define its host type $\text{HOST}(\tau) \in \text{Type}_{\Sigma_{\text{init}}}$ and its relativization predicate on that type, $\text{REL}(\tau) : \text{HOST}(\tau) \Rightarrow \text{bool}$ (where $\text{REL}(\tau) \in \text{Term}_{\Sigma_{\text{init}}}$)
- for each term $t : \tau$, we define its unfolding $\text{UNF}(t) : \text{HOST}(\tau)$ (where $\text{UNF}(t) \in \text{Term}_{\Sigma_{\text{init}}}$)

For instances c_σ of constants $c : \tau$ defined by equations $c_\tau \equiv t$, $\text{UNF}(c_\sigma)$ will be recursively defined as $\text{UNF}(t[\rho])$ where ρ is the substitution that makes σ an instance of τ (i.e., $\sigma \leq_\rho \tau$). In other words, we unfold c_σ with the appropriately substituted equation defining c .

Since UNF is applied to arbitrary terms, not only to constants, we must indicate its recursive behavior for all term constructs. Abstraction and application are handled as expected, but variables raise a subtle issue, with global implications on our overall proof strategy. What should $\text{UNF}(x_\sigma)$ be? $x_{\text{HOST}(\sigma)}$ is an immediate candidate. However, this will not work, since a crucial property that we will need about our translation is that it observes membership to types, in that it maps terms of a given type to terms satisfying that type's representing predicate:

(F1) The relativization predicates hold on translated items, i.e., $\text{REL}(\sigma) \text{ UNF}(t)$ is deducible (in initial HOL) for each term $t : \sigma$.

In particular, $\text{REL}(\sigma) \text{ UNF}(x_\sigma)$ should be deducible. To enforce this, we define $\text{UNF}(x_\sigma)$ to be either $x_{\text{HOST}(\sigma)}$ if $\text{REL}(\sigma) x_{\text{HOST}(\sigma)}$ or else any item for which $\text{REL}(\sigma)$ holds. This is expressible using the if-then-else and Choice operators: $\text{if_t_e}(\text{REL}(\sigma) x_{\text{HOST}(\sigma)}) x_{\text{HOST}(\sigma)} (\varepsilon \text{REL}(\sigma))$. By the Choice axiom, $\text{REL}(\sigma)$ holds for $\varepsilon \text{REL}(\sigma)$ just in case $\text{REL}(\sigma)$ is nonempty. So to achieve the goal of ensuring $\text{REL}(\sigma)$ holds for x_σ , we need:

(F2) The relativization predicates are nonempty, i.e., $\exists x_{\text{HOST}(\sigma)}. \text{REL}(\sigma) x$ is deducible.

Another way to regard this property is as a reflection of the HOL types being nonempty—a faithful relativization should of course follow suit.

Because of the way we apply these definitions recursively to the type and term constructs, the desired Σ_{init} -formula φ' corresponding to φ will be $\text{UNF}(\varphi)$. For example, as one would expect, the

unfoldings of $\forall x_\sigma. \varphi x$ and $\exists x_\sigma. \varphi x$ will be (deduction-equivalent to) $\forall x_{\text{HOST}(\sigma)}. \text{REL}(\sigma) x \longrightarrow \text{UNF}(\varphi) x$ and $\exists x_{\text{HOST}(\sigma)}. \text{REL}(\sigma) x \wedge \text{UNF}(\varphi) x$, respectively. Hence, for us meta-safety over initial HOL will mean:

(MS) For all $\varphi \in \text{Fmla}_\Sigma$, $D \vdash_\Sigma \varphi$ implies $\vdash_{\Sigma_{\text{init}}} \text{UNF}(\varphi)$.

This property is indeed a type-aware version of what Wenzel calls meta-safety: $\text{UNF}(\varphi)$ replaces each defined constant with a term as in Wenzel's concept, and replaces each defined type with a tandem of a host type and a relativization predicate.

To help proving (MS), we will also have lemmas about the good behavior of the translation functions HOST , UNF and REL with respect to the main ingredients of HOL deduction:

(F3) The translation functions preserve variable freshness and commute with substitution.

The order in which we will have to prove these facts has superficially circular dependencies. As discussed, we need (F2) for proving (F1). Moreover, (F1) is needed to prove (F3), more precisely, to make sure that UNF commutes with substitution for the delicate case of variables x_σ . In turn, (F3) is used for (MS). But to prove (F2), the nonemptiness of the relativization predicates, we seem to need (MS). Indeed, for the case of a type τ defined by $\tau \equiv t$ with $t : \sigma \Rightarrow \text{bool}$, the natural choice for $\text{REL}(\tau)$ is the conjunction of $\text{REL}(\sigma)$ and $\text{UNF}(t)$: gathering recursively whatever comes from the potential definition of σ or of its component types and adding the translation of τ 's own defining predicate. So, in an inductive proof of (F2), we will need to deduce $\exists x_{\text{HOST}(\sigma)}. \text{REL}(\sigma) x \wedge \text{UNF}(t) x$. The only fact that can help here is that this formula is (equivalent to) $\text{UNF}(\varphi)$, where φ is $\exists x_\sigma. t x$. Since φ is the non-emptiness claim for the new type τ , it is deducible (according to Def. 2(2)). So we would like to apply (MS) here for obtaining that $\text{UNF}(\varphi)$ is deducible.

In summary, we would need (F2) to prove (MS) and (MS) to prove (F2). The way out of this loop is a *gradual* approach: we will not define a single version of the translation functions, but one version, HOST_i , UNF_i and REL_i , for each subset $\{\text{def}_1, \dots, \text{def}_i\}$ of D with $i \leq n$. This way, we can use (MS) for i to prove (F2) for $i + 1$.

Finally, we must take into account a phenomenon we have ignored so far: the presence of *declarations* in addition to definitions. Let c be a declared constant of type σ . What should its unfolding $\text{UNF}(c_\sigma)$ be? A possibility is to acknowledge c as an irreducible entity, and define $\text{UNF}(c_\sigma) = c_{\text{HOST}(\sigma)}$. However, this way our desirable property (F1), here, $\text{REL}(\sigma) c_{\text{HOST}(\sigma)}$, will not be provable, since nothing prevents the “uninterpreted” items $c_{\text{HOST}(\sigma)}$ from laying outside of the relativization predicate. Another alternative is to define $\text{UNF}(c_\sigma)$ as an arbitrary element satisfying $\text{REL}(\sigma)$, via Choice, i.e., as $\epsilon \text{REL}(\sigma)$. But this would mean that UNF will artificially identify several distinct constants, e.g., $\text{UNF}(c_\sigma) = \text{UNF}(d_\sigma)$ for any two declared constants c_σ and d_σ —besides being unnatural, this situation would become difficult to handle for Isabelle/HOL, since it would introduce a breach in monotony: When declaring c_σ and d_σ , their unfoldings would be equal, but at a later stage one of them could get defined, breaking this equality. (Incidentally, the semantic version of this problem makes Isabelle/HOL constant definitions non-conservative model-theoretically.)

In summary, we wish to preserve the identity of the declared constants $c : \sigma$, while still enforcing $\text{REL}(\sigma) \text{UNF}(c_\sigma)$. We achieve this by treating c_σ in a guarded fashion, similarly to the variables x_σ , i.e., taking $\text{UNF}(c_\sigma)$ to be $\text{if_t_e} (\text{REL}(\sigma) c_{\text{HOST}(\sigma)}) c_{\text{HOST}(\sigma)} (\epsilon \text{REL}(\sigma))$. As for the declared (but not defined) types, these can be kept in the signature without causing any problems. Since we do not eliminate the declared constants and types, in the statement (MS) of meta-safety we must replace Σ_{init} with a suitable signature Δ containing the declared items. Declarations can be

intertwined with definitions, in particular, constants of *defined* types can be declared—so what we need is not only to collect all declared constants, but to also translate their types to the host types.

The rest of this section will unfold the ideas described above. First we illustrate the ideas by some examples, then we formally define and study the translations, culminating with proofs of meta-safety and conservativity.

4.2 Examples

We start with an extensive example that has only definitions, no declarations:

EXAMPLE 4. Let Σ be the extension of the initial signature with:

- the nullary type constructors nat and zfun
- the constants $\text{absnat} : \text{ind} \Rightarrow \text{nat}$, $\text{z} : \text{nat}$ and $\text{repzfun} : \text{zfun} \Rightarrow (\text{nat} \Rightarrow \text{nat})$

Let $D = \{\text{def}_i \mid i \in \{1, \dots, 5\}\}$, such that:

- def_1 is $\text{nat} \equiv t_1$, where $t_1 : \text{ind} \Rightarrow \text{bool}$ is a term in the initial signature (namely, the predicate representing the intersection of all predicates that holds for 0 and are closed under Suc)
- def_2 is $\text{absnat} \equiv t_2$, where t_2 is $\varepsilon t'_2$, with $t'_2 : (\text{ind} \Rightarrow \text{nat}) \Rightarrow \text{bool}$ a predicate (stating that its argument function is a bijection between the elements of ind that satisfy t_1 and nat)
- def_3 is $\text{z} \equiv t_3$, where t_3 is $\text{absnat } 0$
- def_4 is $\text{zfun} \equiv t_4$, where $t_4 : (\text{nat} \Rightarrow \text{nat}) \Rightarrow \text{bool}$ is $\lambda f_{\text{nat} \Rightarrow \text{nat}}. f \text{ z} = \text{z}$
- def_5 is $\text{repzfun} \equiv t_5$, where t_5 is $\varepsilon t'_5$ with $t'_5 : (\text{zfun} \Rightarrow (\text{nat} \Rightarrow \text{nat})) \Rightarrow \text{bool}$ a predicate stating that its argument is one-to-one and its image is included in t_4

Thus, there are no (non-defined but) declared items, and the chain $\Sigma_{\text{init}} = \Sigma_0 \subseteq \Sigma^1 \subseteq \Sigma_1 \subseteq \dots \subseteq \Sigma^5 \subseteq \Sigma_5$ consists of the following signatures, where we omit repeating the arities and the types:

$$\begin{array}{ll} \Sigma^1 = \Sigma_0 = \Sigma_{\text{init}} & \Sigma^4 = \Sigma_3 = \Sigma^3 \cup \{\text{z}\} \\ \Sigma^2 = \Sigma_1 = \Sigma^1 \cup \{\text{nat}\} & \Sigma^5 = \Sigma_4 = \Sigma^4 \cup \{\text{zfun}\} \\ \Sigma^3 = \Sigma_2 = \Sigma^2 \cup \{\text{absnat}\} & \Sigma = \Sigma_5 = \Sigma^5 \cup \{\text{repzfun}\} \end{array}$$

Incidentally, this example shows the standard procedure of bootstrapping natural numbers in HOL: The type nat is defined by carving out, from HOL's built-in infinite type ind , the smallest set closed under zero and successor. Using the Choice operator, we define the abstraction function absnat as a surjection that respects nat 's defining predicate t_1 . (The opposite injection can of course also be defined, but is omitted here.) The version of zero for naturals, $\text{z} : \text{nat}$, is defined by applying the abstraction to the built-in zero from ind .

Subsequently, another type is introduced, zfun , of zero-preserving functions between naturals, defined by carving out from the type $\text{nat} \Rightarrow \text{nat}$ the set of those functions that map z to z . For this type, we define the representation function repzfun to its defining type $\text{nat} \Rightarrow \text{nat}$. Note that the way to apply an element of zfun to a natural is to apply its representation.

We will focus on evaluating $\text{UNF}(f_{\text{zfun}})$, where we write UNF for the last (widest-reaching) unfolding function UNF_5 (and similarly for HOST and REL). As discussed, since f_{zfun} is a variable, $\text{UNF}(f_{\text{zfun}})$ will be a term for which $\text{REL}(\text{zfun})$ is guaranteed to hold (provided the predicate in nonempty): $\text{if_t_e}(\text{REL}(\text{zfun}) f_{\text{HOST}(\text{zfun})}) f_{\text{HOST}(\text{zfun})} (\varepsilon \text{REL}(\text{zfun}))$. Now, looking at the types in definitions def_1 and def_4 , we can compute the host of zfun :

$$\text{HOST}(\text{zfun}) = \text{HOST}(\text{nat} \Rightarrow \text{nat}) = \text{HOST}(\text{nat}) \Rightarrow \text{HOST}(\text{nat}) = \text{ind} \Rightarrow \text{ind}$$

Thus, $\text{REL}(\text{zfun})$ is a predicate on $\text{ind} \Rightarrow \text{ind}$. But what does it say? To evaluate $\text{REL}(\text{zfun})$, we again look at the definitions def_1 and def_4 , this time also factoring in their terms, t_1 and t_4 :

$$\begin{aligned} \text{REL}(\text{zfun}) &= \\ \lambda g_{\text{HOST}(\text{zfun})}. \text{REL}(\text{nat} \Rightarrow \text{nat}) g \wedge \text{UNF}(t_4) g &= \\ \lambda g_{\text{HOST}(\text{zfun})}. (\text{REL}(\text{nat}) \Rightarrow \text{REL}(\text{nat})) g \wedge \text{UNF}(t_4) g &= \\ \lambda g_{\text{ind} \Rightarrow \text{ind}}. (\text{UNF}(t_1) \Rightarrow \text{UNF}(t_1)) g \wedge \text{UNF}(t_4) g &= \\ \lambda g_{\text{ind} \Rightarrow \text{ind}}. (t_1 \Rightarrow t_1) g \wedge \text{UNF}(t_4) g & \end{aligned}$$

where, for a predicate such as $t_1 : \text{ind} \Rightarrow \text{bool}$, $t_1 \Rightarrow t_1$ denotes its lifting to functions:

$$\lambda g_{\text{ind} \Rightarrow \text{ind}}. \forall x_{\text{ind}}. t_1 x \longrightarrow t_1 (g x).$$

(Note that $\text{UNF}(t_1) = t_1$ since t_1 is in the initial signature.) Thus, $\text{REL}(\text{zfun}) g_{\text{ind} \Rightarrow \text{ind}}$ states that g preserves t_1 (the isomorphic image of nat in ind) and that $\text{UNF}(t_4) g$ holds, where t_4 is the isomorphic image of zfun in $\text{nat} \Rightarrow \text{nat}$. This shows how, when evaluating REL , nested type definitions lead to the accumulation of their defining predicates, each lifted if necessary along the encountered function-space structure.

We can prove $D \vdash_{\Sigma} \varphi$, where φ is $\forall f_{\text{zfun}}. \text{repzfun } f_{\text{zfun}} z = z$ with f_{zfun} a variable, i.e., that the items in zfun indeed map zero to zero. By our meta-safety result, we will infer $\vdash_{\Sigma_{\text{init}}} \text{UNF}(\varphi)$, which boils down to a tautology: that all functions from ind to ind that preserve the natural-number-predicate and preserve 0 also preserve 0.

We conclude with an example showing how declarations affect the target signature:

EXAMPLE 5. Consider the following extension of Example 4: After def_4 , a declaration of a constant $c : \text{zfun}$ is performed. Thus, Σ^5 is no longer equal to Σ_4 , but is $\Sigma_4 \cup \{(c, \text{zfun})\}$.

What should be the signature of $\text{UNF}(c_{\text{zfun}})$? Since c has no definition, it will not be compiled away by unfolding. However, we are required to compile away its type zfun , which is a defined type. So it is natural to have $\text{UNF}(c_{\text{zfun}}) = \text{if_t_e } (\text{REL}(\text{zfun}) c_{\text{HOST}(\text{zfun})}) c_{\text{HOST}(\text{zfun})} (\varepsilon \text{REL}(\text{zfun})) = \text{if_t_e } (\text{REL}(\text{zfun}) c_{\text{ind} \Rightarrow \text{ind}}) c_{\text{ind} \Rightarrow \text{ind}} (\varepsilon \text{REL}(\text{zfun}))$. However, none of the existing signatures contains a constant $c : \text{ind} \Rightarrow \text{ind}$.

Consequently, we must create a signature Δ that extends Σ_{init} with all the declared constants but having HOST -translated types, and, similarly, with all the declared type constructors. In general, the translations will target this signature rather than Σ_{init} .

4.3 Formal Definition of the Translations and Meta-Safety

We will write D_i for the current definitional theory at moment i , $\{\text{def}_1, \dots, \text{def}_i\}$. Thus, we have $D = D_n$. As discussed, we will define deduction-preserving translations of the Σ -types and Σ -terms into Δ -types and Δ -terms, where Δ will be a suitable signature that collects all the declared items. We proceed gradually, considering Σ_i one i at a time, eventually reaching $\Sigma = \Sigma_n$.

For each $i \in \{1, \dots, n\}$, we define the signature Δ^i (collecting the declared items from Σ^i with their types translated to their host types), together with the function $\text{HOST}_i : \text{Type}_{\Sigma_i} \Rightarrow \text{Type}_{\Delta^i}$ (producing the host types) as follows:

- Δ^1 is Σ^1
- Δ^{i+1} is Δ^i extended with:
 - all the type constructors $k \in \Sigma^{i+1} \setminus \Sigma_i$
 - for all constants $c \in \Sigma^{i+1} \setminus \Sigma_i$ of type σ , a constant c of type $\text{HOST}_i(\sigma)$
- HOST_i is defined as in Fig. 2, recursively on types

- (H1) $\text{HOST}_i(\alpha) = \alpha$
(H2) $\text{HOST}_i((\sigma_1, \dots, \sigma_m) k) = (\text{HOST}_i(\sigma_1), \dots, \text{HOST}_i(\sigma_m)) k$,
if $k \in \Sigma^1 \cup \bigcup_{i'=2}^i (\Sigma^{i'} \setminus \Sigma_{i'-1})$
(H3) $\text{HOST}_i((\sigma_1, \dots, \sigma_m) k) = \text{HOST}_i(\sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m])$,
if $(\alpha_1, \dots, \alpha_m) k \equiv t$ is in D_i and $t : \sigma \Rightarrow \text{bool}$
- (R1) $\text{REL}_i(\sigma) = \lambda x_\sigma. \text{True}$, if $\sigma \in \text{TVar} \cup \{\text{bool}, \text{inf}\}$
(R2) $\text{REL}_i(\sigma_1 \Rightarrow \sigma_2) = \lambda f_{\text{HOST}_i(\sigma_1) \Rightarrow \text{HOST}_i(\sigma_2)}. \forall x_{\text{HOST}_i(\sigma_1)}. \text{REL}_i(\sigma_1) x \longrightarrow \text{REL}_i(\sigma_2) (f x)$
(R3) $\text{REL}_i((\sigma_1, \dots, \sigma_m) k) = \lambda x_{(\text{HOST}_i(\sigma_1), \dots, \text{HOST}_i(\sigma_m)) k}. \text{True}$, if $k \in \bigcup_{i'=1}^i (\Sigma^{i'} \setminus \Sigma_{i'-1})$
(R4) $\text{REL}_i((\sigma_1, \dots, \sigma_m) k) = \lambda x_{\text{HOST}_i(\sigma')}. \text{REL}_i(\sigma') x \wedge \text{UNF}_i(t') x$,
if $(\alpha_1, \dots, \alpha_m) k \equiv t$ is in D_i and $t : \sigma \Rightarrow \text{bool}$,
where $\sigma' = \sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m]$ and $t' = t[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m]$
- (U1) $\text{UNF}_i(x_\sigma) = \text{if_t_e} (\text{REL}_i(\sigma) x_{\text{HOST}_i(\sigma)}) x (\varepsilon \text{REL}_i(\sigma))$
(U2) $\text{UNF}_i(c_\sigma) = c_{\text{HOST}_i(\sigma)}$, if $c \in \Sigma_{\text{init}}$
(U3) $\text{UNF}_i(c_\sigma) = \text{if_t_e} (\text{REL}_i(\sigma) c_{\text{HOST}_i(\sigma)}) c_{\text{HOST}_i(\sigma)} (\varepsilon \text{REL}_i(\sigma))$, if $c \in \bigcup_{i'=1}^i (\Sigma^{i'} \setminus \Sigma_{i'-1})$
(U4) $\text{UNF}_i(c_\sigma) = \text{UNF}_i(t[\rho])$, if $c_\tau \equiv t$ is in D_i and $\sigma \leq_\rho \tau$
(U5) $\text{UNF}_i(t_1 t_2) = \text{UNF}_i(t_1) \text{UNF}_i(t_2)$
(U6) $\text{UNF}_i(\lambda x_\sigma. t) = \lambda x_{\text{HOST}_i(\sigma)}. \text{UNF}_i(t)$

Fig. 2. Definition of the translation functions

On defined types (i.e., types having a defined type constructor on top, clause (H3)), HOST_i behaves as prescribed in Section 4.1, recursively calling itself for the defining type. Upon encountering built-in or declared type constructors, i.e., belonging to some $\Sigma^{i'}$ for $i' \leq i$, but not to the corresponding $\Sigma_{i'-1}$ (clause (H2)), HOST_i delves into the subexpressions.

Next, *mutually* recursively on Σ_i -types and Σ_i -terms, we define a function returning the relativization predicate of a type, $\text{REL}_i : \text{Type}_{\Sigma_i} \rightarrow \text{Term}_{\Delta_i}$, and one returning the unfolded term, $\text{UNF}_i : \text{Term}_{\Sigma_i} \rightarrow \text{Term}_{\Delta_i}$. Their definition is shown in Fig. 2. Again, they behave as prescribed in Section 4.1. In particular, REL_i is naturally lifted to function spaces (clause (R2)) and accumulates defining predicates, as shown in clause (R4)—here, the substitution $\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m$ stems from an instance of the defined type, $(\alpha_1, \dots, \alpha_m) k$. Type variables and declared types are treated as black boxes, so REL_i is vacuously true for them, just like for the built-in types `bool` and `ind` (clauses (R1) and (R3)). Note that, while (H2) refers to declared or built-in type constructors, (R3) only refers to declared ones—it explicitly excludes Σ_{init} .

As discussed in Section 4.1, UNF_i treats type variables and declared constants in a guarded fashion (clauses (U1) and (U3)), and distributes over application and abstraction (clauses (U5) and (U6)). Moreover, UNF_i merely calls HOST_i for built-in constants (clause (U2)). Finally, UNF_i unfolds the definitions of defined constants, as shown in clause (U4). In that clause, c_τ and $\rho \upharpoonright_{\text{TV}(c_\tau)}$ (the restriction of ρ to $\text{TV}(c_\tau)$) are uniquely determined by c_σ ; and since $\text{TV}(t) \subseteq \text{TV}(c_\sigma)$ (by Def. 1), it follows that $t[\rho]$ is also uniquely determined by c_σ .

Obviously, these functions can reach their purpose only if they are total functions. i.e., their recursive evaluation process terminates for all inputs. This is what we prove in the next subsection.

Assuming totality, we have all the prerequisites to formulate meta-safety. We let UNF be UNF_n , the function that unfolds all definitions in $D = D_n$, and Δ be Δ^n , the signature collecting all the declared items in Σ .

DEFINITION 6. D is said to be a *meta-safe extension of HOL-with-declarations* if, for all $\varphi \in \text{Fmla}_\Delta$, it holds that $D \vdash_\Sigma \varphi$ implies $\vdash_\Delta \text{UNF}(\varphi)$.

4.4 Totality of the Translations

The goal of this subsection is to prove:

PROP 7. The following hold:

- (1) The function HOST_i is total, i.e., its recursive calls terminate.
- (2) The functions REL_i and UNF_i are total, i.e., their mutually recursive calls terminate.

The concepts we use in the proof of this proposition, in particular, the *definitional dependency relation*, will be also relevant in Section 5, when we attend to Isabelle/HOL.

To prove (1), we must show that the call graph of HOST_i , namely, the relation \blacktriangleright_i defined by:

$$\begin{aligned} (\sigma_1, \dots, \sigma_m) k \blacktriangleright_i \sigma_j & \text{ if } k \in \Sigma^i \\ (\sigma_1, \dots, \sigma_m) k \blacktriangleright_i \sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m] & \text{ if } (\alpha_1, \dots, \alpha_m) k \equiv t \text{ is in } D_i \text{ and } t : \sigma \Rightarrow \text{bool} \end{aligned}$$

is terminating. This is easily done by defining a lexicographic order based on the order in which the items were defined, i.e., the indexes of the definitions def_i in which they appear. (Details are given in the appendix.)

To prove (2), we will exhibit a terminating relation \blacktriangleright_i that captures the mutual call graph of REL_i and UNF_i . We take \blacktriangleright_i to be the union $\equiv_i^\downarrow \cup \triangleright$, where \equiv_i^\downarrow and \triangleright are defined below. The relation \triangleright consists of the structurally recursive calls of REL_i and UNF_i , from clauses (R2), (U1), (U5) and (U6):

$$\sigma_1 \Rightarrow \sigma_2 \triangleright \sigma_1 \quad \sigma_1 \Rightarrow \sigma_2 \triangleright \sigma_2 \quad x_\sigma \triangleright \sigma \quad t_1 t_2 \triangleright t_1 \quad t_1 t_2 \triangleright t_2 \quad \lambda x_\sigma. t \triangleright t$$

Moreover, \equiv_i^\downarrow captures the recursive calls corresponding to defined items, from (R4) and (U4). Given $u, v \in \text{Type}_{\Sigma_i} \cup \text{Term}_{\Sigma_i}$, $u \equiv_i^\downarrow v$ states that there exists a definition $u' \equiv v'$ in D_i and a type substitution ρ such that $u = \rho(u')$ and $v = \rho(v')$.

Thus, the totality of REL_i and UNF_i is reduced to the termination of \blacktriangleright_i . In order to prove the latter, we will introduce a more basic relation: the dependency relation between non-built-in items induced by definitions in D_i . We let $\text{Type}_{\Sigma_i}^\bullet$ be the set of Σ_i -types that have a non-built-in type constructor at the top, and $\text{CInst}_{\Sigma_i}^\bullet$ be the set of instances of non-built-in constants. Given any term t , we let $\text{types}^\bullet(t)$ be the set of all types from $\text{Type}_{\Sigma_i}^\bullet$ appearing in t and $\text{cinsts}^\bullet(t)$ be the set of all constant instances from $\text{CInst}_{\Sigma_i}^\bullet$ appearing in t . (The appendix gives the formal definition of these operators.)

DEFINITION 8. The *dependency relation* \rightsquigarrow_i on $\text{Type}_{\Sigma_i}^\bullet \cup \text{CInst}_{\Sigma_i}^\bullet$ is defined as follows: $u \rightsquigarrow_i v$ iff there exists in D_i a definition of the form $u \equiv t$ such that $v \in \text{cinsts}^\bullet(t) \cup \text{types}^\bullet(t)$.

We write $\rightsquigarrow_i^\downarrow$ for the (type-)substitutive closure of \rightsquigarrow_i , defined as follows: $u \rightsquigarrow_i^\downarrow v$ iff there exist u', v' and a type substitution ρ such that $u = u'[\rho]$, $v = v'[\rho]$ and $u' \rightsquigarrow_i v'$. Since HOL with definitions is well-known to be consistent, one would expect that definitions cannot introduce infinite (including cyclic) chains of dependencies. This can indeed be proved by a lexicographic argument, again taking advantage of the definitional order:

LEMMA 9. The relation $\rightsquigarrow_i^\downarrow$ is terminating.

The next observation connects \blacktriangleright_i and $\rightsquigarrow_i^\downarrow$, via \triangleright^* (the transitive closure of \triangleright):

LEMMA 10. If $u, v \in \text{Type}_{\Sigma_i}^\bullet \cup \text{CInst}_{\Sigma_i}^\bullet$ and $u \equiv_i^\downarrow t \triangleright^* v$, then $u \rightsquigarrow_i^\downarrow v$

Now we can reduce the termination of \blacktriangleright_i to that of $\rightsquigarrow_i^\downarrow$, hence prove the former:

LEMMA 11. The relation \blacktriangleright_i is terminating.

This concludes the proof of Prop. 7.

4.5 Basic Properties of the Translations

As envisioned in Section 4.1, the translations are extensions of each other and preserve type membership:

LEMMA 12. Assume $i \leq n - 1$. The following hold:

- (1) If $\sigma \in \text{Type}_{\Sigma_i}$, then $\text{HOST}_{i+1}(\sigma) = \text{HOST}_i(\sigma)$
- (2) If $\sigma \in \text{Type}_{\Sigma_i}$, then $\text{REL}_{i+1}(\sigma) = \text{REL}_i(\sigma)$.
- (3) If $t \in \text{Term}_{\Sigma_i}$, then $\text{UNF}_{i+1}(t) = \text{UNF}_i(t)$.

LEMMA 13. If $\sigma \in \text{Type}_{\Sigma_i}$, $t \in \text{Type}_{\Sigma_i}$ and $t : \sigma$, then $\text{REL}_i(\sigma) : \text{HOST}_i(\sigma) \Rightarrow \text{bool}$ and $\text{UNF}_i(t) : \text{HOST}_i(\sigma)$.

For items in the initial signature, the behavior of the translations is either idle (for HOST_i and UNF_i) or trivial (for REL_i):

LEMMA 14. The following hold:

- (1) If $\sigma \in \text{Type}_{\Sigma_{\text{init}}}$, then $\text{HOST}_i(\sigma) = \sigma$
- (2) If $\sigma \in \text{Type}_{\Sigma_{\text{init}}}$, then $\vdash_{\Sigma_{\text{init}}} \text{REL}_i(\sigma) = \lambda x_{\text{HOST}_i(\sigma)}. \text{True}$
- (3) If $t \in \text{Term}_{\Sigma_{\text{init}}}$ and t is well-typed, then $\vdash_{\Sigma_{\text{init}}} \text{UNF}_i(t) = t$

Other easy, but important properties state that the translations do not introduce new variables or type variables and commute with *type* substitution:

LEMMA 15. The following hold for all $\sigma \in \text{Type}_{\Sigma_i}$ and $t \in \text{Term}_{\Sigma_i}$:

- (1) $\text{TV}(\text{HOST}_i(\sigma)) \subseteq \text{TV}(\sigma)$
- (2) $\text{TV}(\text{REL}_i(\sigma)) \subseteq \text{TV}(\sigma)$ and $\text{FV}(\text{REL}_i(\sigma)) = \emptyset$
- (3) $\text{TV}(\text{UNF}_i(t)) \subseteq \text{TV}(t)$ and $\text{FV}(\text{UNF}_i(t)) = \{x_{\text{HOST}_i(\sigma)} \mid x_\sigma \in \text{FV}(t)\}$

LEMMA 16. The following hold for all $\sigma, \tau \in \text{Type}_{\Sigma_i}$ and $t \in \text{Term}_{\Sigma_i}$:

- (1) $\text{HOST}_i(\sigma[\tau/\alpha]) = \text{HOST}_i(\sigma)[\text{HOST}_i(\tau)/\alpha]$
- (2) $\text{REL}_i(\sigma[\tau/\alpha]) = \text{REL}_i(\sigma)[\text{HOST}_i(\tau)/\alpha]$
- (3) $\text{UNF}_i(t[\tau/\alpha]) = \text{UNF}_i(t)[\text{HOST}_i(\tau)/\alpha]$

4.6 Main Results

We are now ready to finalize the plan set out in Section 4.1. The following facts in Lemma 17 are stated and proved in the delicate order prescribed there. Fact (4) corresponds to part of (F3) (the remaining parts being covered by Lemmas 15 and 16). Moreover, (2) corresponds to (F2), (3) to (F1), and (5) to (MS). Finally, (1) states deducibility of the translated nonemptiness statement, identified in Section 4.1 as an intermediate fact leading from (MS) to (F2).

LEMMA 17. Let $i \in \{1, \dots, n\}$. The following hold for all $\sigma, \tau \in \text{Type}_{\Sigma_i}$, $t, t' \in \text{Term}_{\Sigma_i}$ and $\varphi \in \text{Fmla}_{\Sigma_i}$:

- (1) If $\tau \equiv t$ is a type definition in D_i with $t : \sigma \Rightarrow \text{bool}$, then $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x$
- (2) $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x$
- (3) If $t : \sigma$, then $\vdash_{\Delta^i} \text{REL}_i(\sigma) \text{UNF}_i(t)$
- (4) If $t' : \sigma$, then $\vdash_{\Delta^i} \text{UNF}_i(t[t'/x_\sigma]) = \text{UNF}_i(t)[\text{UNF}_i(t')/x_{\text{HOST}_i(\sigma)}]$

(5) If $D_i \vdash_{\Sigma_i} \varphi$, then $\vdash_{\Delta^i} \text{UNF}_i(\varphi)$

Proof. The facts follow by induction on i . More precisely, let $(j)_i$ denote fact (j) for a given layer i . We prove:

- that $(1)_i$ holds;
- that, for any $i \in \{1, \dots, n\}$:
 - $(1)_i$ implies $(2)_i$ implies $(3)_i$ implies $(4)_i$;
 - $(2)_i$ and $(4)_i$ imply $(5)_i$;
- that, for any $i \in \{1, \dots, n-1\}$, $(5)_i$ implies $(1)_{i+1}$.

We only show proof sketches for the two most crucial of these implications. (The appendix discusses the others.)

$(1)_i$ implies $(2)_i$: Assuming $(1)_i$, we prove $(2)_i$ by structural induction on σ . The only interesting case is when the type is defined, i.e., has a defined type constructor on top (dealt with in clause (R4)). We need to show $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma')}. \text{REL}_i(\sigma') x \wedge \text{UNF}_i(t') x$, where $(\alpha_1, \dots, \alpha_m) k \equiv t$ is in D_i and $t : \sigma \Rightarrow \text{bool}$, $\sigma' = \sigma[(\sigma_j/\alpha_j)_j]$, and $t' = t[(\sigma_j/\alpha_j)_j]$.

By $(1)_i$, we have $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x$. By the type substitution rule (T-INST) applied m times (once for each $\text{HOST}_i(\sigma_j)/\alpha_j$), we have $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)[(\text{HOST}_i(\sigma_j)/\alpha_j)_j]}. \text{REL}_i(\sigma)[(\text{HOST}_i(\sigma_j)/\alpha_j)_j] x \wedge \text{UNF}_i(t)[(\text{HOST}_i(\sigma_j)/\alpha_j)_j] x$. Using Lemma 16 m times (once for each σ_j/α_j), we obtain $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma[(\sigma_j/\alpha_j)_j])}. \text{REL}_i(\sigma[(\sigma_j/\alpha_j)_j]) x \wedge \text{UNF}_i(t[(\sigma_j/\alpha_j)_j]) x$, which implies $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma')}. \text{REL}_i(\sigma') x \wedge \text{UNF}_i(t') x$, as desired.

$(2)_i$ and $(4)_i$ imply $(5)_i$: Assume $(2)_i$ and $(4)_i$. By induction on the definition of HOL deduction (\vdash), we prove a slight generalization of $(5)_i$, namely: We assume $\Gamma \cup \{\varphi\} \subseteq \text{Fmla}_{\Sigma_i}$ and $D_i; \Gamma \vdash_{\Sigma_i} \varphi$, and prove $\emptyset; \text{UNF}_i(\Gamma) \vdash_{\Delta^i} \text{UNF}_i(\varphi)$. We distinguish different cases, according to the last applied rule in inferring $\Gamma \cup \{\varphi\} \subseteq \text{Fmla}_{\Sigma_i}$:

(FACT): We need to prove $\emptyset; \text{UNF}_i(\Gamma) \vdash_{\Delta^i} \text{UNF}_i(\varphi)$, assuming $\varphi \in \text{Ax} \cup D_i$. First, assume $\varphi \in D$. Then $\varphi = u \equiv t \in D_i$. We have two subcases:

(A) u is a constant c_σ . Then $\text{UNF}_i(\varphi)$ is the formula $\text{UNF}_i(c_\sigma) = \text{UNF}_i(t)$. And since $\text{UNF}_i(c_\sigma)$ and $\text{UNF}_i(t)$ are (syntactically) equal, the desired fact follows by the HOL reflexivity rule.

(B) u is a type τ of the form $(\alpha_1, \dots, \alpha_m) k$ and $t : \sigma \Rightarrow \text{bool}$. Then, by the definition of UNF_i and of the \forall and \exists constructs, $\text{UNF}_i(\varphi)$ is deduction-equivalent to the formula

$$\begin{aligned} & \exists \text{rep}_{\text{HOST}_i(\sigma) \Rightarrow \text{HOST}_i(\sigma)}. (\forall x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x \longrightarrow \text{REL}_i(\sigma) (\text{rep } x)) \\ & \wedge \\ & \forall x_{\text{HOST}_i(\sigma)}, y_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x \wedge \text{REL}_i(\sigma) y \wedge \text{UNF}_i(t) y \wedge \text{rep } x = \text{rep } y \longrightarrow x = y \\ & \wedge \\ & \forall y_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) y \longrightarrow (\text{UNF}_i(t) y \leftrightarrow (\exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x \wedge y = \text{rep } x)) \end{aligned}$$

where the first conjunct comes from the relativization of $\tau \Rightarrow \sigma$, the second from unfolding $\text{One_One}_{\text{rep}}$, and the third from unfolding $\forall y_{\sigma}. t y \leftrightarrow (\exists x_{\tau}. y = \text{rep } x)$ (in Def. 1). This states the following (in a verbose fashion): There exists $\text{rep} : \text{HOST}_i(\sigma) \Rightarrow \text{HOST}_i(\sigma)$ which is one-to-one on the intersection of $\text{REL}_i(\sigma)$ and $\text{UNF}_i(t)$ and the image of this intersection through rep is the intersection itself. This is of course deducible in HOL, taking rep as the identity function.

Now, assume $\varphi \in \text{Ax}$. Then $\varphi \in \text{Fmla}_{\Sigma_{\text{init}}}$, hence, by Lemma 14(3), $\vdash_{\Delta^i} \text{UNF}_i(\varphi) = \varphi$. And since also $\emptyset; \text{UNF}_i(\Gamma) \vdash_{\Delta^i} \varphi$ is true by (FACT), the desired fact follows using the HOL equality rules.

(ASSUM): Follows by applying (ASSUM).

(T-INST): Courtesy of UNF_i commuting with type substitution (Lemma 16(3)) and preserving freshness (Lemma 15(3)).

(INST): Courtesy of UNF_i commuting with substitution (point (4)_{*i*}) and preserving freshness (Lemma 15(3)).

(BETA), (EXT), (IMPI) and (MP): Courtesy of UNF_i commuting with substitution, preserving freshness, and distributing (by definition) over abstractions, applications and implications. \square

As a particular case of this lemma's point (5), we have:

THEOREM 18. D is a meta-safe extension of HOL-with-declarations.

Thus, we can compile away all the definitions of D , which leaves us with types and terms over the signature Δ containing declarations only. With the definitions out of our way, it remains to show that *declarations* are conservative, which is much easier:

LEMMA 19. If $\varphi \in \text{Fmla}_{\Sigma_{\text{init}}}$ and $\vdash_{\Delta} \varphi$, then $\vdash_{\Sigma_{\text{init}}} \varphi$.

Proof. Assume $\vdash_{\Delta} \varphi$. In the proof tree for this fact, we replace:

- (1) all occurrences of any declared constant instance c_{σ} by a fresh variable x_{σ}
- (2) all occurrences of any declared type constructor k of arity m by a built-in type expression of arity n , e.g., $(\sigma_1, \dots, \sigma_m)k$ is replaced by $\sigma_1 \Rightarrow \dots \Rightarrow \sigma_m$

Then the resulted proof tree constitutes a proof of $\vdash_{\Sigma_{\text{init}}} \varphi$. \square

Finally, we can prove overall conservativity:

THEOREM 20. D is a conservative extension of initial HOL.

Proof. Assume $D \vdash_{\Sigma} \varphi$, where $\varphi \in \text{Fmla}_{\Sigma_{\text{init}}}$. By Theorem 18, we have $\vdash_{\Delta} \text{UNF}(\varphi)$. Moreover, by Lemma 14(3), we have $\vdash_{\Sigma_{\text{init}}} \text{UNF}(\varphi) = \varphi$, hence, a fortiori, $\vdash_{\Delta} \text{UNF}(\varphi) = \varphi$. From these two, we obtain $\vdash_{\Delta} \varphi$. With Lemma 19, we obtain $\vdash_{\Sigma_{\text{init}}} \varphi$, as desired. \square

4.7 Abstract Constant Definition Mechanisms

As definitional schemes for constants, we have only looked into the traditional *equational* ones, implemented in most HOL provers. Two non-equational schemes have also been designed [Arthan 2014], and are available in HOL4, HOL Light and ProofPower-HOL: “new specification” and “gen new specification.” They allow for more abstract (under)specification of constants.

However, these schemes have been shown not to increase expressiveness: “new specification” can be over-approximated by traditional definitions and the use of the Choice operator, and “gen new specification” is an admissible rule in HOL with “new specification” [Arthan 2014; Kumar et al. 2014]. Hence our results cater for them.

5 CONSERVATIVITY OF ISABELLE/HOL DEFINITIONS

As mentioned in the introduction, Isabelle/HOL allows more flexible constant definitions than HOL, in that it enables ad hoc overloaded definitions. For example, one can declare a polymorphic constant, such as $\leq : \alpha \Rightarrow \alpha \text{ bool}$, and at later times (perhaps after some other type and constant definitions and declarations have been performed) define different, non-overlapping instances of it: \leq_{nat} as the standard order on natural numbers, \leq_{bool} as implication, etc. Even recursive overloading is allowed, e.g., one can define $\leq_{\alpha \text{ list}}$ as the component-wise extension of \leq_{α} to α list:

$$xs \leq_{\alpha \text{ list}} ys \equiv \text{length } xs = \text{length } ys \wedge (\forall i < \text{length } xs. xs_i \leq_{\alpha} ys_i)$$

This means that now constant definitions no longer require the constant to be *fresh*. In fact, we are no longer speaking of constant definitions, but of constant *instance* definitions: The above examples do not define the overall constant \leq , but various instances of it, \leq_{nat} , \leq_{bool} and \leq_{list} .

DEFINITION 21. Given a non-built-in constant c , a type $\sigma \leq \text{tpOf}(c)$ and a closed term $t : \sigma$, we let $c_\sigma \equiv t$ denote the formula $c_\sigma = t$. We call $c_\sigma \equiv t$ a *constant-instance definition* provided $\text{TV}(t) \subseteq \text{TV}(c_\sigma)$.

To compensate for the lack of freshness from constant-instance definitions, the Isabelle/HOL system performs some global syntactic checks, making sure that defined instances do not overlap (i.e., definitions are *orthogonal*) and that the dependency relation \rightsquigarrow_n from Def. 8, terminates [Kunčar 2015; Kunčar and Popescu 2015, 2017]. (Recall that $D = D_n$, hence \rightsquigarrow_n is the dependency induced by D , i.e., by all the considered definitions.) Formally:

DEFINITION 22. An *Isabelle/HOL-well-formed definitional theory* is a set D of type and constant-instance definitions over Σ such that:

- It satisfies all the conditions of Def. 2, except that it is *not* required that, in condition (1.2), c be fresh, i.e., it is *not* required that $c \notin \Sigma^i$
- It is orthogonal: For all constants c , if c_σ and c_τ appear in two definitions in D , then $\sigma \# \tau$
- Its induced dependency relation \rightsquigarrow_n is terminating

We wish to prove meta-safety and conservativity results similar to the ones for traditional HOL. To this end, we fix an Isabelle/HOL-well-formed definitional theory D and look into the results of Section 4 to see what can be reused—as it turns out, quite a lot.

First, the (type-translated) declaration signatures Δ^i and the translation functions HOST_i , REL_i and UNF_i are defined in the same way. The orthogonality assumption in Def. 22 ensures that, in clause (U4) from the definition of UNF_i , the choice of t is unique (whereas before, this was simply ensured by c appearing on the left in at most one definition). The notion of meta-safety is then defined in the same way. Thanks to \rightsquigarrow_n being terminating, all the dependency relations \rightsquigarrow_i , which are included in \rightsquigarrow_n , are also terminating. Then all the results in Section 4.4 hold, leading to the totality of the translation functions. Furthermore, almost all the lemmas in Section 4.5 go through undisturbed, because they do not need the freshness assumption $c \notin \Sigma^i$.

The only losses are parts of Lemmas 12 (extension of the translations from i to $i + 1$) and 16 (commutation with type substitution), namely, points (2) and (3) of these lemmas—which deal with REL_i and UNF_i . We first look at Lemma 16.

While HOST_i still commutes with substitution, this is no longer the case for REL_i and UNF_i . Essentially, $\text{UNF}_i(\sigma[\tau/\alpha]) = \text{UNF}_i(\sigma)[\text{HOST}_i(\tau)/\alpha]$ now fails because $\text{UNF}_i(\sigma[\tau/\alpha])$ gets to unfold more constant-instance definitions than $\text{UNF}_i(\sigma)$. So the difference is that, for the constant instances $c_{\sigma'}$ occurring in σ that happen to have a definition of one of their instances, say, $c_{\sigma''} \equiv t$ with $\sigma'' \leq \sigma'$, activated by the substitution τ/α (meaning we have $\sigma'[\tau/\alpha] \leq \sigma''$, but $\sigma' \not\leq \sigma''$), $\text{UNF}_i(\sigma[\tau/\alpha])$ will unfold $c_{\sigma'}$ into the corresponding instance of $\text{UNF}(t)$, whereas $\text{UNF}_i(\sigma)[\text{HOST}_i(\tau)/\alpha]$ will replace $c_{\sigma'}$ with $\text{if_t_e}(\text{REL}_i(\sigma') c_{\text{HOST}_i(\sigma')} c_{\text{HOST}_i(\sigma')} (\varepsilon \text{REL}_i(\sigma')))$. (And since REL_i depends recursively on UNF_i , the former will also fail to commute with type substitution.)

EXAMPLE 23. To Example 4's signature, we add a declared constant c of polymorphic type α and a definition of its nat-instance, $c_{\text{nat}} \equiv z$. We have $\text{UNF}(c_\alpha[\text{nat}/\alpha]) = \text{UNF}(c_{\text{nat}}) = \text{UNF}(z)$, whereas $\text{UNF}(c_\alpha)[\text{HOST}(\text{nat})/\alpha] = (\text{if_t_e}(\text{REL}(\alpha) c_{\text{HOST}(\alpha)}) c_{\text{HOST}(\alpha)} (\varepsilon \text{REL}(\alpha))) [\text{ind}/\alpha] = (\text{if_t_e} \text{True } c_\alpha (\varepsilon (\lambda x. \text{True}))) [\text{ind}/\alpha] =_{\text{HOL}} c_\alpha[\text{ind}/\alpha] = c_{\text{ind}}$, where we wrote $=_{\text{HOL}}$ for HOL-provable equality (in the current signature). We do not need to evaluate $\text{UNF}(z)$ in order to see that it cannot be equal, not even HOL-provably equal, to c_{ind} . Indeed, the constant c was not even present in the signature when z was defined, so $\text{UNF}(z)$ cannot be connected to c_{ind} .

Fortunately, we can amend this mismatch “after the fact” by replacing $c_{\text{HOST}_i(\sigma')}$ with $\text{UNF}_i(c_{\sigma''})$ in $\text{UNF}_i(\sigma)[\text{HOST}_i(\tau)/\alpha]$ for all instances $c_{\sigma''}$ (with $\sigma'' \leq \sigma'$) of all defined constant instances $c_{\sigma'}$.

In the above example, this means replacing c_{ind} with $\text{UNF}(c_{\text{nat}})$, i.e., with $\text{UNF}(z)$. To express this formally, we define a *constant-instance substitution* to be a function $\gamma : \text{CInst}_{\Delta^i}^{\bullet} \Rightarrow \text{Term}_{\Delta^i}$ such that, for all $c_{\sigma} \in \text{CInst}_{\Delta^i}^{\bullet}$, $\gamma(c_{\sigma})$ is a closed term and $\text{TV}(\gamma(c)) \subseteq \text{TV}(c)$ —thus assigning a term to any instance of a non-built-in, i.e., declared constant in Δ^i . Using a notation similar to variable substitution, we write $\sigma[[\gamma]]$ and $t[[\gamma]]$ for the effect of performing γ everywhere inside the type σ or the term t .

LEMMA 24. There exists a constant-instance substitution γ such that:

- (1) $\vdash_{\Delta^i} \text{REL}_i(\sigma[\tau/\alpha]) = \text{REL}_i(\sigma)[\text{HOST}_i(\tau)/\alpha] [[\gamma]]$
- (2) $\vdash_{\Delta^i} \text{UNF}_i(t[\tau/\alpha]) = \text{UNF}_i(t)[\text{HOST}_i(\tau)/\alpha] [[\gamma]]$

Now, the question is whether the partial conflation offered by Lemma 24, a quasi-commutativity property for REL_i and UNF_i , can replace full commutativity towards the central goal in Lemma 17, namely, point (5) (which ensures meta-safety). Answering this will require some proof mining.

The only usage of Lemma 16 was for (1)_{*i*} implies (2)_{*i*} (which is part of an implication chain leading to (4)_{*i*}; and both (2)_{*i*} and (4)_{*i*} are used for (5)_{*i*}). There, we used Lemma 16 m times to infer $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma')}. \text{REL}_i(\sigma') x \wedge \text{UNF}_i(t') x$ from $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x$. So we actually need a weaker statement:

LEMMA 25. If $\vdash_{\Delta^i} \text{UNF}_i(\varphi)$, then $\vdash_{\Delta^i} \text{UNF}_i(\varphi[\sigma/\alpha])$.

Proof. By Lemma 24(2), we have a constant-instance substitution γ such that $\text{UNF}_i(\varphi[\sigma/\alpha]) = \text{UNF}_i(\varphi)[\text{HOST}_i(\sigma)/\alpha] [[\gamma]]$. And since $\vdash_{\Delta^i} \text{UNF}_i(\varphi)[\text{HOST}_i(\sigma)/\alpha]$ follows from $\vdash_{\Delta^i} \text{UNF}_i(\varphi)$ by the type substitution rule (T-INST), it would suffice to have the following: For all constant-instance substitutions γ and Δ^i -formulas φ , $\vdash_{\Delta^i} \varphi$ implies $\vdash_{\Delta^i} \varphi[[\gamma]]$. In words, if we substitute some (undefined) constant instances with terms of the same type we do not lose provability. This follows by routine induction on the definition of deduction. \square

For Lemma 12, the situation is quite similar to that of Lemma 16. This time, it is not substitution that can enable additional unfoldings, but a newly added instance definition $c_{\sigma} \equiv t$ at layer $i + 1$ for a constant c that already existed at layer i . Moreover, when we look at how we employed Lemma 12 in the proof of our main chain of results in Lemma 17, we discover a similar pattern: We only use that UNF_{i+1} and REL_{i+1} extend UNF_i and REL_i in the proof of (5)_{*i*} implies (1)_{*i+1*}, where we needed that deduction at layer $i + 1$ is implied by deduction at layer i . By a similar trick as before, this can be proved using a weaker quasi-commutativity property.

LEMMA 26. If $\varphi \in \text{Fmla}_{\Sigma^i}$, and $\vdash_{\Delta^i} \text{UNF}_i(\varphi)$, then $\vdash_{\Delta^{i+1}} \text{UNF}_{i+1}(\varphi)$.

Proof. If def_{i+1} is a type definition, then UNF_{i+1} and REL_{i+1} do extend UNF_i and REL_i , so the desired fact follows trivially. Now, assume def_{i+1} is a constant-instance definition $c_{\sigma} \equiv t$. Similarly to the proof of Lemma 24(2), we obtain a constant-instance substitution γ such that $\text{UNF}_{i+1}(\varphi) = \text{UNF}_i(\varphi) [[\gamma]]$, namely, γ maps each $d_{\text{HOST}(\tau[\rho])}$ to $\text{UNF}_{i+1}(s[\rho])$ where $d_{\tau} \equiv s$ are the constant definitions in D_{i+1} . (We need to do this replacement to all defined constant instances, not just c_{σ} , since other definitions from D_i may have already relied on c_{σ} .) And since, as we have seen, constant-instance substitution preserves deduction, we obtain our desired fact. \square

Lemma 25 and 26 reflect a concession made to Isabelle/HOL's ad hoc overloading: We can no longer exhibit a precise structural relationship between $\text{UNF}_i(\varphi)$ on the one hand and $\text{UNF}_i(\varphi[\sigma/\alpha])$ or $\text{UNF}_{i+1}(\varphi)$ on the other, but we can prove that the latter are “at least as deducible as the former.” This would not have been possible had we not treated declared constants in a guarded fashion in the UNF_i clause (U3) (see the discussion on page 11).

Thus, we were able to recover Lemma 17’s point (5), leading to meta-safety. And since the other ingredients in the proof of Theorem 20 are also available (including Lemma 19, which is independent of the definitional mechanisms), we infer conservativity. We obtained:

THEOREM 27. Theorems 18 and 20 still hold if we assume that D is an Isabelle/HOL-well-formed definitional theory.

6 CONCLUDING REMARKS

We have resolved an open problem, relevant for the foundation of HOL-based theorem provers, including our favorite one, Isabelle/HOL: We showed that the definitional mechanisms in such provers are meta-safe and conservative over pure HOL, i.e., are truly “definitional.” Our result has for HOL a foundational status analogous to strong normalization results for type theory.

Our statement of meta-safety is calibrated to what we believe is the key desirable property: that definitions can *all* be compiled away, without loss of provability. An even more general statement would involve compiling away *some* definitions $E \subseteq D$ only, and translating any statement involving all definitions into one involving all definitions but those in E .

However, even the formulation of meta-safety seems problematic here: Say we define the polymorphic type αk as the following subset of bool:

$$\text{if_t_e (cardinal } \alpha = 3) \{ \text{True, False} \} \{ \text{True} \}$$

Then we define the type l as the subset $\{1, 2, 3\}$ of ind . Stating that l ’s definition is meta-safe over k ’s definition would require us, e.g., to find a host type for $l k$ without being allowed to unfold k . The only sensible choice for the host would be $\text{ind } k$, which is not suitable since $l k$ is larger than $\text{ind } k$: The former has two elements, whereas the latter has one. This means that we cannot relativize $l k$ as a predicate on $\text{ind } k$. Abstractly, the problem is that we cannot lift relativization predicates from the types with which αk may be instantiated (such as l). If each HOL type constructor had the structure of a relator (endofunctor on the category of sets having relations as morphisms), the lifting would be possible in a canonical way. And most useful types in HOL, e.g., all combinations of inductive and coinductive datatypes and function spaces, are in fact relators [Traytel et al. 2012]. However, `typedef` *can* introduce (rather strange looking) non-relators: k is an example of a type constructor that cannot be organized as a relator.

Also, if k were merely *declared*, we would not have a problem, since then we could treat it as a black box that renders $\text{ind } k$ and $l k$ indistinguishable; so we could take the latter’s relativization predicate to be vacuously true. In our meta-safety theorems, we employed this trick to cover declarations intermixed with definitions.

Notwithstanding the difficulty with formulating a more general meta-safety, we believe *conservativity* holds more generally, but requires a different proof technique.

REFERENCES

- Andreas Abel, Thierry Coquand, and Peter Dybjer. 2007. Normalization by Evaluation for Martin-Lof Type Theory with Typed Equality Judgements. In *LICS*. 3–12.
- Mark Adams. 2010. Introducing HOL Zero (Extended Abstract). In *ICMS ’10*. Springer.
- Thorsten Altenkirch. 1993. Proving Strong Normalization of CC by Modifying Realizability Semantics. In *TYPES*. 3–18.
- Rob Arthan. 2014. “HOL Constant Definition Done Right”. In *ITP*. 531–536.
- R. D. Arthan. 2004. Some Mathematical Case Studies in ProofPower–HOL. In *TPHOLS*.
- Andrea Asperti, Wilmer Ricciotti, Claudio Sacerdoti Coen, and Enrico Tassi. 2011. The Matita Interactive Theorem Prover. In *CADE*. 64–69.
- Bruno Barras. 2010. Sets in Coq, Coq in Sets. *Journal of Formalized Reasoning* 3, 1 (2010).
- Yves Bertot and Pierre Casteran. 2004. *Interactive Theorem Proving and Program Development. Coq’Art: The Calculus of Inductive Constructions*. Springer.

- Jasmin Christian Blanchette, Johannes Hölzl, Andreas Lochbihler, Lorenz Panny, Andrei Popescu, and Dmitriy Traytel. 2014. Truly Modular (Co)datatypes for Isabelle/HOL. In *ITP*, Vol. 8558. 93–110.
- Ana Bove, Peter Dybjer, and Ulf Norell. A Brief Overview of Agda—A Functional Language with Dependent Types. In *TPHOLS 2009*.
- Alonzo Church. 1940. A Formulation of the Simple Theory of Types. *The Journal of Symbolic Logic* 5, 2 (1940), 56–68.
- Thierry Coquand, Jean Gallier, and Le Chesnay Cedex. 1990. A Proof of Strong Normalization For the Theory of Constructions Using a Kripke-Like Interpretation. In *Workshop on Logical Frameworks*.
- Thierry Coquand and Arnaud Spiwack. 2006. A Proof of Strong Normalisation using Domain Theory. In *LICS*. 307–316.
- Javier Esparza, Peter Lammich, René Neumann, Tobias Nipkow, Alexander Schimpf, and Jan-Georg Smaus. 2013. A Fully Verified Executable LTL Model Checker. In *CAV*. 463–478.
- J.H. Geuvers. 1993. *Logics and Type systems*. Ph.D. Dissertation. University of Nijmegen.
- M. J. C. Gordon and T. F. Melham (Eds.). 1993. *Introduction to HOL: A Theorem Proving Environment for Higher Order Logic*. Cambridge University Press.
- Florian Haftmann and Makarius Wenzel. 2006. Constructive Type Classes in Isabelle.. In *TYPES*.
- John Harrison. 1996. HOL Light: A Tutorial Introduction. In *FMCAD '96*. Springer.
- John Harrison. 2006. Towards self-verification of HOL Light. In *IJCAR 2006*. Springer.
- Leon Henkin. 1949. The Completeness of the First-Order Functional Calculus. *J. Symbolic Logic* 14, 3 (09 1949), 159–166.
- Isabelle. 2016. The Isabelle Library. (2016). <https://isabelle.in.tum.de/dist/library/HOL/index.html>.
- Matt Kaufmann, Panagiotis Manolios, and J Strother Moore. 2000. *Computer-Aided Reasoning: An Approach*. Kluwer Academic Publishers.
- Gerwin Klein, June Andronick, Kevin Elphinstone, Gernot Heiser, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. 2010. seL4: formal verification of an operating-system kernel. *Commun. ACM* 53, 6 (2010), 107–115.
- Gerwin Klein, Tobias Nipkow, Larry Paulson, and René Thiemann (eds.). 2016. Isabelle’s Archive of Formal Proofs. (2016).
- Alexander Krauss. 2009. *Automating recursive definitions and termination proofs in higher-order logic*. Ph.D. Dissertation. Technical University Munich.
- Ramana Kumar, Rob Arthan, Magnus O. Myreen, and Scott Owens. 2014. HOL with Definitions: Semantics, Soundness, and a Verified Implementation. In *ITP*.
- Ramana Kumar, Rob Arthan, Magnus O. Myreen, and Scott Owens. 2016. Self-Formalisation of Higher-Order Logic - Semantics, Soundness, and a Verified Implementation. *J. Autom. Reasoning* 56, 3 (2016), 221–259.
- Ondřej Kunčar. 2015. Correctness of Isabelle’s Cyclicity Checker: Implementability of Overloading in Proof Assistants. In *CPP*. 85–94.
- Ondřej Kunčar and Andrei Popescu. 2015. A Consistent Foundation for Isabelle/HOL. In *ITP*. 234–252.
- Ondřej Kunčar and Andrei Popescu. 2017. Comprehending Isabelle/HOL’s Consistency. In *ESOP*. To appear. Preprint available at http://andreipopescu.uk/pdf/compr_IsabelleHOL_cons.pdf.
- Andreas Lochbihler. 2010. Verifying a Compiler for Java Threads. In *ESOP*. 427–447.
- Thomas F. Melham. 1989. Automating Recursive Type Definitions in Higher Order Logic. In *Current Trends in Hardware Verification and Automated Theorem Proving*. 341–386.
- Magnus O. Myreen and Jared Davis. 2014. The Reflective Milawa Theorem Prover Is Sound - (Down to the Machine Code That Runs It). In *ITP*. 421–436.
- Tobias Nipkow, Lawrence Paulson, and Markus Wenzel. 2002. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*. LNCS, Vol. 2283. Springer.
- Tobias Nipkow and Gregor Snelting. 1991. Type Classes and Overloading Resolution via Order-Sorted Unification. In *Functional Programming Languages and Computer Architecture*.
- Steven Obua. 2006. Checking Conservativity of Overloaded Definitions in Higher-Order Logic.. In *RTA*.
- Sam Owre and Natarajan Shankar. 1999. The Formal Semantics of PVS. (1999). SRI technical report. <http://www.csl.sri.com/papers/csl-97-2/>.
- Lawrence C. Paulson. 1990. A formulation of the simple theory of types (for Isabelle). In *COLOG-88*. 246–274.
- Lawrence C. Paulson. 2010. Three Years of Experience with Sledgehammer, a Practical Link between Automatic and Interactive Theorem Provers. In *Proceedings of the 2nd Workshop on Practical Aspects of Automated Reasoning, PAAR-2010, Edinburgh, Scotland, UK, July 14, 2010*. 1–10.
- Lawrence C. Paulson. 2015. A Mechanised Proof of Gödel’s Incompleteness Theorems Using Nominal Isabelle. *J. Autom. Reasoning* 55, 1 (2015), 1–37.
- A. Pitts. 1993. *Introduction to HOL: A Theorem Proving Environment for Higher Order Logic*, Chapter The HOL Logic, 191–232. In Gordon and Melham [Gordon and Melham 1993].

- Donald Sannella and Andrzej Tarlecki. 2012. *Foundations of Algebraic Specification and Formal Software Development*. Springer. I–XVI, 1–581 pages.
- Konrad Slind and Michael Norrish. 2008. "A Brief Overview of HOL4". In *TPHOLS*. 28–32.
- Dmitriy Traytel, Andrei Popescu, and Jasmin Christian Blanchette. 2012. Foundational, Compositional (Co)datatypes for Higher-Order Logic: Category Theory Applied to Theorem Proving. In *LICS*. 596–605.
- D. A. Turner. 2004. Total Functional Programming. *J. UCS* 10, 7 (2004), 751–768.
- Markus Wenzel. 1997. Type Classes and Overloading in Higher-Order Logic.. In *TPHOLS*.
- Markus Wenzel. 1999. Isar - A Generic Interpretative Approach to Readable Formal Proof Documents. In *TPHOLS*. 167–184.
- Makarius Wenzel. 2014. System description: Isabelle/jEdit in 2014. In *UTP*. 84–94.
- Freek Wiedijk. 2009. Stateless HOL. In *TYPES*. 47–61.
- Burkhart Wolff. 2015. Isabelle Foundation & Certification. (2015). Archived at <https://lists.cam.ac.uk/pipermail/cl-isabelle-users/2015-September/thread.html>.

A MORE DETAILS ON HOL

It is well-known (and easy to prove) that substitution respects typing:

LEMMA 28. If $t : \sigma$, then $t[\rho] : \sigma[\rho]$.

When writing concrete terms or formulas, we take the following conventions:

- We omit redundantly indicating the types of the variables, e.g., we shall write $\lambda x_\sigma. x$ instead of $\lambda x_\sigma. x_\sigma$.
- We omit redundantly indicating the types of the variables and constants in terms if they can be inferred by typing rules, e.g., we shall write $\lambda x. (y_{\sigma \Rightarrow \tau} x)$ instead of $\lambda x_\sigma. (y_{\sigma \Rightarrow \tau} x)$ or $\varepsilon(\lambda x_\sigma. P x)$ instead of $\varepsilon_{(\sigma \Rightarrow \text{bool}) \Rightarrow \sigma}(\lambda x_\sigma. P_{\sigma \Rightarrow \text{bool}} x)$.
- We write $\lambda x_\sigma y_\tau. t$ instead of $\lambda x_\sigma. \lambda y_\tau. t$
- We apply the constants \longrightarrow and $=$ in an infix manner, e.g., we shall write $t_\sigma = s$ instead of $= t_\sigma s$.

The formula connectives and quantifiers are defined as abbreviations in the usual way, starting from the implication and equality primitives:

$$\begin{aligned} \text{True} &= (\lambda x_{\text{bool}}. x) = (\lambda x_{\text{bool}}. x) \\ \text{and} &= \lambda p_{\text{bool}} q_{\text{bool}}. \text{All} (\lambda f_{\text{bool} \Rightarrow \text{bool} \Rightarrow \text{bool}}. f p q) = (\lambda f. f \text{True True}) \\ \text{All} &= \lambda p_{\alpha \Rightarrow \text{bool}}. p = (\lambda x. \text{True}) \\ \text{Implies} &= \lambda p_{\text{bool}} q_{\text{bool}}. \text{and } p q = p \\ &(\text{in what follows, we write } p \longrightarrow q \text{ instead of } \text{Implies } p q) \\ \text{Ex} &= \lambda p_{\alpha \Rightarrow \text{bool}}. \text{All} (\lambda q. (\text{All} (\lambda x. p x \longrightarrow q)) \longrightarrow q) \\ \text{False} &= \text{All} (\lambda p_{\text{bool}}. p) \\ \text{not} &= \lambda p. p \longrightarrow \text{False} \\ \text{or} &= \lambda p q. \text{All} (\lambda r. (p \longrightarrow r) \longrightarrow ((q \longrightarrow r) \longrightarrow r)) \end{aligned}$$

It is easy to see that the above terms are closed and well-typed as follows:

- True, False : bool
- not : bool \Rightarrow bool
- and, or : bool \Rightarrow bool \Rightarrow bool
- All, Ex : ($\alpha \Rightarrow$ bool) \Rightarrow bool

As customary, we write:

- $\forall x_\alpha. t$ instead of All $(\lambda x_\alpha. t)$
- $\exists x_\alpha. t$ instead of Ex $(\lambda x_\alpha. t)$
- $\neg \varphi$ instead of not φ
- $\varphi \wedge \chi$ instead of and $\varphi \chi$
- $\varphi \vee \chi$ instead of or $\varphi \chi$

The HOL axioms, forming the set Ax, are the following:

- Equality Axioms:
 - refl: $x_\alpha = x$
 - subst: $x_\alpha = y \longrightarrow p x \longrightarrow p y$
- Infinity Axioms:
 - suc_inj: $\text{suc } x = \text{suc } y \longrightarrow x = y$
 - suc_not_zero: $\neg \text{suc } x = \text{zero}$
- Choice:

some_intro: $p_{\alpha \Rightarrow \text{bool}} x \longrightarrow p (\varepsilon p)$

Above, refl and subst axiomatize equality. suc_inj and suc_not_zero ensure that ind is an infinite type. some_intro regulates the behavior of the Hilbert Choice operator. The principle of excluded middle, $(b = \text{True}) \vee (b = \text{False})$, follows from the axiom of choice—this makes HOL a classical logic.

B DETAILED DEFINITION OF THE OPERATORS USED IN THE DEPENDENCY RELATION

Note that the types^\bullet operator is overloaded for types and terms.

$$\begin{aligned}
 \text{types}^\bullet(\alpha) &= \{\alpha\} & \text{types}^\bullet(x_\sigma) &= \text{types}^\bullet(\sigma) \\
 \text{types}^\bullet(\text{bool}) &= \emptyset & \text{types}^\bullet(c_\sigma) &= \text{types}^\bullet(\sigma) \\
 \text{types}^\bullet(\text{ind}) &= \emptyset & \text{types}^\bullet(t_1 t_2) &= \text{types}^\bullet(t_1) \cup \text{types}^\bullet(t_2) \\
 \text{types}^\bullet(\sigma_1 \Rightarrow \sigma_2) &= \text{types}^\bullet(\sigma_1) \cup \text{types}^\bullet(\sigma_2) & \text{types}^\bullet(\lambda x_\sigma. t) &= \text{types}^\bullet(\sigma) \cup \text{types}^\bullet(t) \\
 \text{types}^\bullet(\bar{\sigma} k) &= \{\bar{\sigma} k\}, \text{ if } k \neq \Rightarrow, \text{ bool, ind} & &
 \end{aligned}$$

$$\begin{aligned}
 \text{cinsts}^\bullet(x_\sigma) &= \emptyset \\
 \text{cinsts}^\bullet(c_\sigma) &= \begin{cases} \{c_\sigma\} & \text{if } c_\sigma \in \text{CInst}^\bullet \\ \emptyset & \text{otherwise} \end{cases} \\
 \text{cinsts}^\bullet(t_1 t_2) &= \text{cinsts}^\bullet(t_1) \cup \text{cinsts}^\bullet(t_2) \\
 \text{cinsts}^\bullet(\lambda x_\sigma. t) &= \text{cinsts}^\bullet(t)
 \end{aligned}$$

C PROOF SKETCHES

In the proofs, we use several induction schemas, fit for the purpose:

- *Well-founded induction* on types and/or terms with respect to one of the (known to be terminating) relations \blacktriangleright_i and \blacktriangleright'_i : Given u , we can assume the property holds for all items u' such that $u \blacktriangleright_i u'$ (or $u \blacktriangleright'_i u'$) and need to prove it for u . So whenever we indicate a proof by well-founded induction, we will implicitly refer to one of these two, namely, to the first when proving something about HOST_i and to the second when proving something about REL_i and/or UNF_i .
- *Structural induction* on types and/or terms: Given u , we can assume the property holds for all immediate subtypes/subterms of u and need to prove it for u .
- *Rule induction* with respect to the definition of typing or the definition of HOL deduction: To conclude that typing or deduction implies a property, we prove that the property is closed under the rules defining typing or deduction.

In all these schemas, (IH) denotes the induction hypothesis.

Proof of point (1) of Prop. 7. We first define, for any type constructor $k \in \Sigma_i$, the operator $\text{depth}_k : \text{Type}_{\Sigma_i} \Rightarrow \mathbb{N}$ to return, for any type, the length of the longest nesting of k 's appearing in it, namely:

$$\begin{aligned}
 \text{depth}_k(\alpha) &= 0 \\
 \text{depth}_k((\sigma_1, \dots, \sigma_m) l) &= \begin{cases} 1 + \max\{\text{depth}_k(\sigma_i) \mid i \in \{1, \dots, m\}\} & \text{if } l = k \\ \max\{\text{depth}_k(\sigma_i) \mid i \in \{1, \dots, m\}\} & \text{if } l \neq k \end{cases}
 \end{aligned}$$

Let K^i and K_i be the sets of type constructors of $\Sigma_1 \cup \bigcup_{i'=1}^i \Sigma^{i'} \setminus \Sigma_{i'-1}$ and Σ_i , respectively. Note that $K^i \subseteq K_i$ and that $K_i \setminus K^i$ contains the defined type constructors, whereas K^i contains the

declared and built-in ones (up to moment i). We chose an arbitrary total order $>$ on K^i , and then extend it to a homonymous total order on K_i , as follows:

- If $k \in K^i$ and $l \in K_i$, then $l > k$
- If $k_1, k_2 \in K_i$, then $k_1 > k_2$ if k_1 was introduced later than k_2 , i.e., if the unique $j_1 \leq i$ such that k_1 appears in the lefthand side of def_{j_1} is greater than the unique $j_2 \leq i$ such that k_2 appears on the lefthand side of def_{j_2}

Since K_i is finite, it has the form $\{k_1, \dots, k_p\}$ with $k_1 > \dots > k_p$. We define the measure $\text{meas} : \text{Type}_{\Sigma_i} \rightarrow \mathbb{N}^p$ by $\text{meas}(\sigma) = (\text{depth}_{k_1}, \dots, \text{depth}_{k_p})$. Finally, we note that \blacktriangleright_i decreases this measure w.r.t. the lexicographic order on \mathbb{N}^p (which ensures its termination). Indeed, we consider the two cases in the definition of \blacktriangleright_i :

- In the first case (given by recursive clause (U3)), all depth_{k_j} remain the same or decrease, and depth_k decreases by 1
- In the second case (given by recursive clause (U4)), we know from the well-foundedness of D that σ only contains type constructors l with $k > l$. Therefore, we have:

$$\begin{aligned} \text{depth}_k((\sigma_1, \dots, \sigma_m) k) &= 1 + \max \{ \text{depth}_k(\sigma_j) \mid j \in \{1, \dots, m\} \} > \\ \max \{ \text{depth}_l(\sigma_j) \mid j \in \{1, \dots, m\} \wedge \alpha_i \in \text{TV}(\sigma) \} &= \text{depth}_k(\sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m]) \end{aligned}$$

Moreover, for any l such that $l > k$, we have:

$$\begin{aligned} \text{depth}_l((\sigma_1, \dots, \sigma_m) k) &= \max \{ \text{depth}_l(\sigma_j) \mid j \in \{1, \dots, m\} \} \geq \\ \max \{ \text{depth}_l(\sigma_j) \mid j \in \{1, \dots, m\} \wedge \alpha_i \in \text{TV}(\sigma) \} &= \text{depth}_l(\sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m]) \end{aligned}$$

Thus, depth_k decreases strictly and, for $l > k$, depth_l remains the same or decreases; this ensures that meas decreases. \square

Proof of Lemma 9. We proceed similarly to the proof of termination for the call graph of HOST_i , but considering Σ_i -constants in addition to Σ_i -type constructors. Similarly to there, for each $e \in K_i \cup \text{Const}_i$, we define $\text{depth}_e : \text{Type}_{\Sigma_i} \cup \text{Term}_{\Sigma_i} \rightarrow \mathbb{N}$, the u -depth of a type or term, to the length of the longest nesting of u 's appearing in it. We similarly order the items in $K_i \cup \text{Const}_i$ by a relation $>$ asking that all defined items are greater than all non-defined ones and a later defined item is greater than an earlier defined one. Assuming $K_i \cup \text{Const}_i$ has the form $\{e_1, \dots, e_p\}$ with $e_1 > \dots > e_p$, we define the measure $\text{meas} : \text{Type}_{\Sigma_i} \cup \text{Term}_{\Sigma_i} \rightarrow \mathbb{N}^p$ by $\text{meas}(v) = (\text{depth}_{e_1}(v), \dots, \text{depth}_{e_p}(v))$.

We show that meas decreases with $\rightsquigarrow_i^\downarrow$ w.r.t. the lexicographic order on \mathbb{N}^p (which makes $\rightsquigarrow_i^\downarrow$ terminating). Assume $u \rightsquigarrow_i^\downarrow v$. Then there exists u', v', ρ such that $u = u'[\rho]$, $v = v'[\rho]$ and $u' \rightsquigarrow_i v'$. Then u' is either a type of the form $(\alpha_1, \dots, \alpha_m)k$ with $k \in K_i$, or a constant instance c_σ ; meaning u is either $(\rho(\alpha_1), \dots, \rho(\alpha_m))k$ or $c_{\sigma[\rho]}$. We let e denote either k or c . In both cases, we have $v' \in \text{types}^\bullet(t) \cup \text{cinsts}^\bullet(t)$ for some $t \in \text{Term}_{\Sigma_i}$. Hence $v \in \text{types}^\bullet(t[\rho]) \cup \text{cinsts}^\bullet(t[\rho])$. By the well-foundedness of D , e is greater than all the type constructors and constants in t (w.r.t. $>$). Then $\text{depth}_e(u) > \text{depth}_e(v)$ and, for all e' such that $e' > e$, $\text{depth}_{e'}(u) \geq \text{depth}_{e'}(v)$. This ensures $\text{meas}(u) > \text{meas}(v)$. \square

Proof of Lemma 10. By routine structural induction on t . \square

Proof of Lemma 11. Let us assume by absurd that \blacktriangleright_i does not terminate. Then there exists an infinite sequence $(w_p)_{p \in \mathbb{N}}$ such that $w_p \blacktriangleright_i w_{p+1}$ for all p . Since \blacktriangleright_i is defined as $\equiv_i^\downarrow \cup \triangleright$ and \triangleright clearly terminates, there must exist an infinite subsequence $(w_{p_j})_{j \in \mathbb{N}}$ such that $w_{p_j} \equiv_i^\downarrow w_{p_{j+1}} \triangleright^* w_{p_{j+1}}$ for

all j . Since from the definition of \equiv_i^\downarrow we have $w_{p_j} \in \text{Type}_{\Sigma_i}^\bullet \cup \text{CInst}_{\Sigma_i}^\bullet$, we obtain from Lemma 10 that $w_{p_j} \rightsquigarrow_i^\downarrow w_{p_{j+1}}$ for all p . This contradicts the termination of $\rightsquigarrow_i^\downarrow$. \square

Proof of Lemma 12. (1): By an easy well-founded induction on σ w.r.t. \blacktriangleright_i , distinguishing between the different cases in the definition of HOST_i and HOST_{i+1} . The definitions are identical for the two functions, and for the defined type case (clause (H3)), we know that k is in Σ_i , ensuring that $(\sigma_1, \dots, \sigma_m)k \equiv t$ is in D_i .

(2) and (3): Similar to (1), by an easy well-founded induction on σ and t w.r.t. \blacktriangleright_i . \square

Proof of Lemma 13. By well-founded induction on σ and t , distinguishing between the different cases in the definitions of REL_i and UNF_i . The proof is routine. We only show the two slightly less obvious cases, where we employ the local notations used in the definitions (e.g., σ', t'):

The defined type case for REL_i (clause (R4)): We know that $(\sigma_1, \dots, \sigma_m)k \blacktriangleright_i \sigma'$ and also that $(\sigma_1, \dots, \sigma_m)k \blacktriangleright_i t'$. Moreover, from $t : \sigma$ we obtain $t' : \sigma'$. Hence, by (IH), we have $\text{REL}_i(\sigma') : \text{HOST}(\sigma') \Rightarrow \text{bool}$ and $\text{UNF}_i(t') : \text{HOST}(\sigma')$. From this, the definition of REL_i and the HOL typing rules, we obtain $\text{REL}_i((\sigma_1, \dots, \sigma_m)k) : \text{HOST}_i(\sigma') \Rightarrow \text{bool}$. Finally, from the definition of HOST_i we have $\text{HOST}_i(\sigma') = \text{HOST}_i((\sigma_1, \dots, \sigma_m)k)$, hence $\text{REL}_i((\sigma_1, \dots, \sigma_m)k) : \text{HOST}_i((\sigma_1, \dots, \sigma_m)k) \Rightarrow \text{bool}$, as desired.

The defined constant case for UNF_i (clause (U4)): We know that $c_\sigma \blacktriangleright_i t[\rho]$. Moreover, since $\sigma = \tau[\rho]$ and $t : \tau$, by Lemma 28 we have that $t[\rho] : \sigma$. By (IH), we have $\text{UNF}(t[\rho]) : \text{HOST}(\sigma)$; and since $\text{UNF}_i(c_\sigma) = \text{UNF}_i(t[\rho])$, we obtain $\text{UNF}_i(c_\sigma) : \text{HOST}(\sigma)$, as desired. \square

Proof of Lemma 14. (1) and (2): Immediate by structural induction on σ .

(3): Immediate by structural induction on t . For the variable case, we use point (2) to obtain $\vdash_{\Sigma_0} \text{UNF}_i(x_\sigma) = x_\sigma$ from the behavior of if-then-else.

(Since Σ_{init} has no defined or declared items, the recursive cases that deal with such items do not occur when applying the translations, and in particular REL_i does not depend recursively of UNF_i . This is why structural induction does the job, so there is no need for the more powerful well-founded induction.) \square

Proof of Lemma 15. Immediate well-founded induction, using the property that definitions do not introduce free term variables or type variables. \square

Proof of Lemma 16. By routine well-founded induction, using the properties of type substitution. For example: In the defined type cases for HOST_i and REL_i (clauses (H3) and (R4)), we use that, if $\text{TV}(\sigma) \subseteq \{\alpha_1, \dots, \alpha_m\}$ (as it is guaranteed by Def. 1), $\sigma[\sigma_1/\alpha_1, \dots, \sigma_m/\alpha_m][\tau/\alpha] = \sigma[(\sigma_1[\tau/\alpha])/\alpha_1, \dots, (\sigma_m[\tau/\alpha])/\alpha_m]$; in the defined constant case for UNF_i (clause (U4)), we use that $t[\rho][\tau/\alpha] = t[\rho \cdot (\tau/\alpha)]$. (Recall that \cdot is the composition of substitutions.) \square

Remaining cases in the proof of Lemma 17.

(1)₁: By the well-formedness of D (Def. 2), we have that $t \in \text{Term}_{\Delta^i}$ and $\sigma \in \text{Term}_{\Delta^i}$, hence $\text{HOST}_0(\sigma) = \sigma$, $\vdash_{\Delta^i} \text{REL}_0(\sigma) = \lambda x_\sigma. \text{True}$ and $\vdash_{\Delta^i} \text{UNF}_0(t) = t$. From this, we obtain that the fact to be proved is equivalent to $\vdash_{\Delta^i} \exists x_\sigma. t x$, which is again true by the well-formedness of D .

Next, we fix $i \in \{1, \dots, n\}$.

(2) _{i} implies (3) _{i} : Assume (2) _{i} . Then (3) _{i} follows by rule induction on the definition of typing. For

the variable case, we use $(2)_i$ and the Choice axiom, which ensure us that $\vdash_{\Delta^i} \text{REL}_i(\sigma)(\varepsilon \text{REL}_i(\sigma))$ holds, hence $\vdash_{\Delta^i} \text{REL}_i(\sigma)(\text{UNF}_i(x_\sigma))$ holds.

$(3)_i$ implies $(4)_i$: Assume $(3)_i$. Then $(4)_i$ follows by well-founded induction on t . The only interesting case is in the variable case (clause (U1)), when the variable coincides with the to-be substituted variable x_σ . Thus, $t = x_\tau$. Here, we need to show $\vdash_{\Delta^i} \text{UNF}_i(t') = \text{if_t_e}(\text{REL}_i(\sigma) \text{UNF}_i(t'))(\text{UNF}_i(t'))(\varepsilon \text{REL}_i(\sigma))$. This follows from the fact that, thanks to $(3)_i$ and $t' : \sigma$, we have $\vdash_{\Delta^i} \text{REL}_i(\sigma) \text{UNF}_i(t')$.

Next, we fix $i \in \{1, \dots, n-1\}$.

$(5)_i$ implies $(1)_{i+1}$: Assume $(5)_i$ and let σ, t be as in the formulation of $(1)_{i+1}$, namely, $\text{def}_{i+1} = \sigma \equiv t$. By the well-formedness of D (Def. 2), we have $D_i \vdash_{\Sigma_i} \exists x_\sigma. t x$. Applying $(5)_i$, we obtain $\vdash_{\Delta^i} \text{UNF}_i(\exists x_\sigma. t x)$. By the definition of the \exists quantifier and the definition of UNF_i , the above is equivalent to $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x_{\text{HOST}_i(\sigma)} \wedge \text{UNF}_i(t) t'$, where t' is $\text{if_t_e}(\text{REL}_i(\sigma) x_{\text{HOST}_i(\sigma)}) x(\varepsilon \text{REL}_i(\sigma))$. By the definition of the if-then-else operator, we can replace t' by x . So the above is further equivalent to $\vdash_{\Delta^i} \exists x_{\text{HOST}_i(\sigma)}. \text{REL}_i(\sigma) x \wedge \text{UNF}_i(t) x$. By Lemma 12 and the fact that $\Delta^i \subseteq \Delta^{i+1}$, the above implies $\vdash_{\Delta^{i+1}} \exists x_{\text{HOST}_{i+1}(\sigma)}. \text{REL}_{i+1}(\sigma) x \wedge \text{UNF}_{i+1}(t) x$, as desired. \square

Proof of Lemma 24. We will write $t_1 =_{\Delta^i} t_2$ instead of $\vdash_{\Delta^i} t_1 = t_2$. We define γ to map each $c_{\text{HOST}_i(\sigma)}$ to $\text{UNF}_i(c_\sigma)$. Thanks to Lemma 15(3), γ is indeed a constant-instance substitution. Now, points (1) and (2) follow by well-founded induction on the mutually recursive definitions of REL_i and UNF_i . The only interesting case is that of defined constants (clause (U4) for UNF_i). Assume $\sigma[\tau/\alpha] = \sigma'[\rho]$, such that $c_{\sigma'} \equiv t \in D_i$. We have two cases:

First, assume $\sigma \leq \sigma'$, say, $\sigma = \sigma'[\rho']$ for some ρ' . Then ρ and $\rho' \cdot (\tau/\alpha)$ are equal on $\text{TV}(\sigma')$, a fortiori, on $\text{TV}(t)$. Hence $t[\rho] = t[\rho' \cdot (\tau/\alpha)]$. i.e., $t[\rho] = t[\rho'][\tau/\alpha]$ (*). Both $\text{UNF}_i(c_\sigma[\tau/\alpha])$ and $\text{UNF}_i(c_\sigma)$ will unfold the definitions of their corresponding instances of c , allowing us to infer the desired fact from the induction hypothesis:

$$\begin{aligned} \text{UNF}_i(c_\sigma[\tau/\alpha]) &= \text{UNF}_i(c_{\sigma[\tau/\alpha]}) = \text{UNF}_i(c_{\sigma'[\rho]}) = (\text{by (U4)}) = \\ \text{UNF}_i(t[\rho]) &= (\text{by (*)}) = \text{UNF}_i(t[\rho'][\tau/\alpha]) =_{\Delta^i} \\ &(\text{by the induction hypothesis}) \\ \text{UNF}_i(t[\rho'])[\text{HOST}_i(\tau)/\alpha][[\gamma]] &= (\text{by (U4)}) = \\ \text{UNF}_i(c_{\sigma'[\rho']})[\text{HOST}_i(\tau)/\alpha][[\gamma]] &= \text{UNF}_i(c_\sigma)[\text{HOST}_i(\tau)/\alpha][[\gamma]] \end{aligned}$$

Next, assume $\sigma \not\leq \sigma'$. Then only $\text{UNF}_i(c_\sigma[\tau/\alpha])$ unfolds the definition of c'_σ , but γ repairs the mismatch. To ease readability, we will write $_[\bullet]$ instead of $_[\text{HOST}_i(\tau)/\alpha]$.

$$\begin{aligned} \text{UNF}_i(c_\sigma[\tau/\alpha]) &=_{\Delta^i} \\ &(\text{since } \vdash_{\Delta^i} \text{REL}_i(\sigma[\tau/\alpha]) \text{UNF}_i(c_{\sigma[\tau/\alpha]}) \text{ holds}) \\ \text{if_t_e}(\text{REL}_i(\sigma[\tau/\alpha]) \text{UNF}_i(c_{\sigma[\tau/\alpha]})) \text{UNF}_i(c_{\sigma[\tau/\alpha]}) &(\varepsilon \text{REL}_i(\sigma[\tau/\alpha])) = \\ &(\text{by the definition of } \gamma) \\ \text{if_t_e}(\text{REL}_i(\sigma[\tau/\alpha]) \gamma(c_{\text{HOST}_i(\sigma[\tau/\alpha])})) \gamma(c_{\text{HOST}_i(\sigma[\tau/\alpha])}) &(\varepsilon \text{REL}_i(\sigma[\tau/\alpha])) = \\ &(\text{by the definition of constant substitution}) \\ \text{if_t_e}(\text{REL}_i(\sigma[\tau/\alpha]) c_{\text{HOST}_i(\sigma[\tau/\alpha])}[[\gamma]]) c_{\text{HOST}_i(\sigma[\tau/\alpha])}[[\gamma]] &(\varepsilon \text{REL}_i(\sigma[\tau/\alpha])) = \\ &(\text{by (IH) for } \text{REL}_i \text{ and } \text{HOST}_i \text{'s commutation with substitution}) \\ \text{if_t_e}(\text{REL}_i(\sigma)[\bullet][[\gamma]] c_{\text{HOST}_i(\sigma)}[\bullet][[\gamma]]) c_{\text{HOST}_i(\sigma)}[\bullet][[\gamma]] &(\varepsilon[\bullet][[\gamma]] \text{REL}_i(\sigma)[\bullet][[\gamma]]) = \\ &(\text{by the definition of constant substitution}) \\ (\text{if_t_e}(\text{REL}_i(\sigma) c_{\text{HOST}_i(\sigma)}) c_{\text{HOST}_i(\sigma)} &(\varepsilon \text{REL}_i(\sigma))) [\bullet][[\gamma]] = \\ &(\text{by (U3)}) \\ \text{UNF}_i(c_\sigma)[\bullet][[\gamma]] & \end{aligned}$$

\square