



Positive Free Higher-Order Logic and its Automation via a Semantical Embedding

Irina Makarenko and Christoph Benzmüller

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

June 15, 2020

Positive Free Higher-Order Logic and its Automation via a Semantical Embedding

Irina Makarenko and Christoph Benzmüller

Freie Universität Berlin, Berlin, Germany
rna.mkr|c.benzmueller@gmail.com

Abstract. Free logics are a family of logics that are free of any existential assumptions. Unlike traditional classical and non-classical logics, they support an elegant modeling of nonexistent objects and partial functions as relevant for a wide range of applications in computer science, philosophy, mathematics, and natural language semantics. While free first-order logic has been addressed in the literature, free higher-order logic has not been studied thoroughly so far. Contributions of this paper include (i) the development of a notion and definition of free higher-order logic in terms of a positive semantics (partly inspired by Farmer’s partial functions version of Church’s simple type theory), (ii) the provision of a faithful shallow semantical embedding of positive free higher-order logic into classical higher-order logic, (iii) the implementation of this embedding in the Isabelle/HOL proof-assistant, and (iv) the exemplary application of our novel reasoning framework for an automated assessment of Prior’s paradox in positive free quantified propositional logics, i.e., a fragment of positive free higher-order logic.

Keywords: Knowledge representation and reasoning, Interactive and automated theorem proving, Philosophical foundations of AI, Partiality and undefinedness, Prior’s paradox

1 Introduction

The proper handling of nonexistence and partiality constitutes a key challenge not only for applications of formal methods in philosophy and mathematics, but also for computational approaches to artificial intelligence and natural language [13, 15, 14]. In a “free logic” terms do not necessarily have to denote existing objects, allowing for theories involving both partial and total functions. For that reason, we expect free logical theories to become increasingly relevant for a wide range of application domains, including but not limited to knowledge representation in artificial intelligence. Free higher-order logics provide elegant solutions to the handling of some well known paradoxes in knowledge representation and reasoning, many of which are beyond first-order logic. Moreover, free logics are well suited to represent abstract objects and to support hypothetical reasoning with fictive (and concrete) entities, and can therefore also be applied in metaphysics, ethics and law.

Modern interactive and automated theorem provers, however, are typically developed for (classical) notions of logic, in which only total functions are

supported natively. Instead of investing time and effort in the development of new theorem provers for free first-order and higher-order logics, a promising approach for the implementation of such logics in existing higher-order theorem provers are *shallow semantical embeddings* (SSEs) [5]. We have successfully applied the SSE approach and we are currently integrating the results reported in this paper in the LogiKEY framework [9] for expressive, pluralistic normative reasoning.

The contributions of this paper manifold. They include (i) the development of a notion and definition of free higher-order logic in terms of a positive semantics (partly inspired by Farmer’s partial functions version of Church’s simple type theory [12]), (ii) the provision of a faithful shallow semantical embedding of positive free higher-order logic into classical higher-order logic, (iii) the implementation of this embedding in the Isabelle/HOL proof-assistant, and (iv) the exemplary application of our novel reasoning framework for an automated assessment of Prior’s paradox [27] in positive free quantified propositional logics, i.e., a fragment of positive free higher-order logic.

Prior, coinciding with Kaplan [17], showed that paradoxes can arise quickly in particular philosophical theories that include both sets and propositions. Bacon, Hawthorne, and Uzquiano [3] discovered that universal instantiation, or, better, the rejection of it, is key to blocking certain paradoxes inherent in such higher-order logics. Logics without existential assumptions, i.e., free logics, just naturally reject the principle of universal instantiation. The family of paradoxes considered by Bacon et al. is represented by what we will call Prior’s paradox/theorem in this paper. Prior’s theorem states:

$$Q \forall p. (Q p \rightarrow \neg p) \rightarrow \exists p. (Q p \wedge p) \wedge \exists p. (Q p \wedge \neg p)$$

Reading $Q p$ as, for instance, ‘Kaplan believes that p ’, Prior’s theorem says that if Kaplan believes that everything that he believes is false, then he believes something true and also believes something false – a logical self-contradiction that should be resolved, and in fact is resolved in free higher-order logic, as we will discuss and demonstrate later in this paper.

The paper structure is as follows: §2 briefly recaps *classical higher-order logic* (HOL), before *positive free higher-order logic* (PFHOL) is introduced in §3. §4 presents a faithful embedding of PFHOL in HOL, and §5 discusses its encoding in Isabelle/HOL. §6 applies the encoded embedding to “solve” Prior’s paradox, and the paper is then concluded in §7.

2 Classical Higher-Order Logic (HOL)

The *simple theory of types* is a classical higher-order logic defined on top of the simply typed λ -calculus. Church’s original definitions as generalized by Henkin [16] to *extensional type theory*, the logical basis of most automated theorem proving systems for higher-order logic, are summarized below.

2.1 Syntax

The main components of Church's type theory are types and terms; more precisely, typed terms. The *set of simple types* \mathcal{T} is freely generated from a set of two base types $\{o, i\}$ and the right associative function type constructor \rightarrow . Intuitively, o is the type of standard truth values and i is the type of individuals. \mathcal{T} is thus defined by $\alpha, \beta := o \mid i \mid (\alpha \rightarrow \beta)$. $\mathcal{T}_o \subset \mathcal{T}$, the set of simple types of (goal) type o , is given by $\beta := o \mid (\alpha \rightarrow \beta)$ (with $\alpha \in \mathcal{T}$). $\mathcal{T}_i \subset \mathcal{T}$, the set of simple types of (goal) type i , is analogously given by $\beta := i \mid (\alpha \rightarrow \beta)$ (with $\alpha \in \mathcal{T}$).

Starting with some nonempty, countable sets of typed constant symbols C_α and some nonempty, countable sets of typed variable symbols V_α , the *simply typed terms* of HOL are defined by the following formation rules (where $\alpha, \beta \in \mathcal{T}$, $P_\alpha \in C_\alpha$ and $x_\alpha \in V_\alpha$):

$$s, t := P_\alpha \mid x_\alpha \mid (s_{\alpha \rightarrow \beta} t_\alpha)_\beta \mid (\lambda x_\alpha. s_\beta)_{\alpha \rightarrow \beta}$$

Moreover, we assume the following constant symbols to be part of our “signature”:
 $\neg_{o \rightarrow o} \in C_{o \rightarrow o}$, $\vee_{o \rightarrow o \rightarrow o} \in C_{o \rightarrow o \rightarrow o}$, $=_{\alpha \rightarrow \alpha \rightarrow o} \in C_{\alpha \rightarrow \alpha \rightarrow o}$, $\forall_{(\alpha \rightarrow o) \rightarrow o} \in C_{(\alpha \rightarrow o) \rightarrow o}$,
 $\iota_{(\alpha \rightarrow o) \rightarrow \alpha} \in C_{(\alpha \rightarrow o) \rightarrow \alpha}$. These constant symbols, which we call logical constants, have a fixed denotation according to their intuitive meaning. For example, the definite description $\iota_{(\alpha \rightarrow o) \rightarrow \alpha} (\lambda x_\alpha. s_o)_{\alpha \rightarrow o}$ denotes the unique object x of type $\alpha \in \mathcal{T}$ satisfying s_o if it exists and an arbitrary object of type α otherwise (we may view this object as an error object). It offers the possibility to define an if-then-else operator as follows:

$$ite_{o \rightarrow \alpha \rightarrow \alpha \rightarrow \alpha} := \lambda s_o. \lambda x_\alpha. \lambda y_\alpha. \iota(\lambda z_\alpha. (s \rightarrow z = x) \wedge (\neg s \rightarrow z = y))$$

Further logical constants can be introduced as abbreviations, e.g., $\wedge_{o \rightarrow o \rightarrow o} := \lambda x_o. \lambda y_o. \neg(\neg x \vee \neg y)$ and $\exists_{(\alpha \rightarrow o) \rightarrow o} := \lambda p_{\alpha \rightarrow o}. \neg \forall(\lambda x_\alpha. \neg(p x))$ with $\alpha \in \mathcal{T}$. All constant symbols, which are not logical constants, are uninterpreted.

Terms of type o are *formulas*, nonformula terms of type $\alpha \in \mathcal{T}_o$ are called *predicates*. Formulas whose leftmost nonparenthesis symbol is either equality or some nonlogical constant or variable are called *atomic formulas*. A variable x is *bound* in a term s if it occurs in the scope of a quantifier in s . x is *free* in s when it is not bound in s .

Some syntax conventions: Type information may be omitted if clear from the context. For each binary operator op with prefix notation $((op s) t)$ we may fall back to its infix notation $(s opt t)$ to improve readability. Likewise, the binder notation $\{\forall, \iota\}(x. s)$ may be used as shorthand for $\{\forall, \iota\}(\lambda x. s)$. In the remainder of this thesis, a matching pair of parentheses in a type or term may be dropped when they are not necessary, assuming that, in addition to the generally known rules, (st) , the application, and $(\lambda x. s)$, function abstraction, are left and right associative, respectively, and that application has a smaller scope than abstraction.

2.2 Semantics

A *frame* D is a set $\{D_\alpha : \alpha \in \mathcal{T}\}$ of nonempty sets (formally *domains*) D_α such that D_i is chosen freely, $D_o = \{T, F\}$ where $T \neq F$ and T represents truth and F

represents falsehood, and $D_{\alpha \rightarrow \beta}$ is the set of all total functions from domain D_α to codomain D_β . A *standard model* is a tuple $M = \langle D, I \rangle$ where D is a frame and I is a family of typed interpretation functions, i.e., $I = \{I_\alpha : \alpha \in \mathcal{T}\}$. Each *interpretation function* I_α maps constants of type α to appropriate objects of D_α . The logical constants $=, \neg, \vee, \forall$ and ι are interpreted as follows:

$$\begin{aligned}
I(=_{\alpha \rightarrow \alpha \rightarrow o}) &:= id \in D_{\alpha \rightarrow \alpha \rightarrow o} \quad \text{s.t. for all } d, d' \in D_\alpha, \\
&\quad id(d, d') = \text{T iff } d \text{ is identical to } d' \\
I(\neg_{o \rightarrow o}) &:= not \in D_{o \rightarrow o} \quad \text{s.t. } not(\text{T}) = \text{F} \text{ and } not(\text{F}) = \text{T} \\
I(\vee_{o \rightarrow o \rightarrow o}) &:= or \in D_{o \rightarrow o \rightarrow o} \quad \text{s.t. } or(v_1, v_2) = \text{T iff } v_1 = \text{T} \text{ or } v_2 = \text{T} \\
I(\forall_{(\alpha \rightarrow o) \rightarrow o}) &:= all \in D_{(\alpha \rightarrow o) \rightarrow o} \quad \text{s.t. for all } f \in D_{\alpha \rightarrow o}, \\
&\quad all(f) = \text{T iff } f(d) = \text{T for all } d \in D_\alpha \\
I(\iota_{(\alpha \rightarrow o) \rightarrow \alpha}) &:= desc \in D_{(\alpha \rightarrow o) \rightarrow \alpha} \quad \text{s.t. for all } f \in D_{\alpha \rightarrow o}, \\
&\quad desc(f) = d \in D_\alpha \text{ if } f(d) = \text{T} \text{ and for} \\
&\quad \text{all } d' \in D_\alpha: \text{ if } f(d') = \text{T}, \text{ then } d' = d, \\
&\quad \text{otherwise } desc(f) = e, \\
&\quad \text{where } e \text{ is an arbitrary object in } D_\alpha
\end{aligned}$$

g_α is a *variable assignment* mapping variables of type α to corresponding objects in D_α . Thus, $g = \{g_\alpha : \alpha \in \mathcal{T}\}$ is a family of typed variable assignments. $g[x \rightarrow d]$ denotes the assignment that is identical to g , except for variable x_α , which is now mapped to d_α . The *value* $\llbracket s_\alpha \rrbracket^{M,g}$ of a HOL term s_α in a standard model M under variable assignment g is an object $d \in D_\alpha$ which is defined as follows:

$$\begin{aligned}
\llbracket P_\alpha \rrbracket^{M,g} &:= I(P_\alpha) \\
\llbracket x_\alpha \rrbracket^{M,g} &:= g(x_\alpha) \\
\llbracket (s_{\alpha \rightarrow \beta} t_\alpha)_\beta \rrbracket^{M,g} &:= \llbracket s_{\alpha \rightarrow \beta} \rrbracket^{M,g} (\llbracket t_\alpha \rrbracket^{M,g}) \\
\llbracket (\lambda x_\alpha. s_\beta)_{\alpha \rightarrow \beta} \rrbracket^{M,g} &:= \text{the function } f \text{ from } D_\alpha \text{ into } D_\beta \\
&\quad \text{s.t. for all } d \in D_\alpha: f(d) = \llbracket s_\beta \rrbracket^{M,g[x \rightarrow d]}
\end{aligned}$$

A formula s is *true* in a standard model M under variable assignment g if and only if $\llbracket s_o \rrbracket^{M,g} = \text{T}$, denoted by $M, g \models s$. A formula s is *valid in* M , denoted by $M \models s$, if and only if $M, g \models s$ for all variable assignments g . Moreover, a formula s is (*generally*) *valid*, denoted by $\models s_o$, if and only if s is valid in all standard models M .

As a consequence of Gödel's incompleteness theorem, Church's type theory with respect to the ordinary semantics based on standard models is incomplete. However, Henkin [16] introduced a generalized notion of model in which the function domains contain enough but not necessarily all functions: In a standard model a domain $D_{\alpha \rightarrow \beta}$ is defined as the set of all total functions from D_α to D_β . In a *Henkin model* (or *general model*) the domains $D_{\alpha \rightarrow \beta}$ in the underlying frame are some nonempty sets of total functions, $D_{\alpha \rightarrow \beta} \subseteq \{f \mid f : D_\alpha \rightarrow D_\beta\}$, containing at least sufficiently many of them so that the valuation function remains total.

For Henkin’s generalized notion of semantics, sound and complete proof calculi exist [16, 1, 2]. Any standard model is obviously also a Henkin model. Hence, any formula that is valid in all Henkin models must be valid in all standard models as well. Therefore, the semantics employed in this paper are Henkin’s general models. For truth, validity, and general validity in a Henkin model, the above definitions are adapted in the obvious way.

For further details on the semantics of HOL, we refer to the literature [7, 6].

3 Positive Free Higher-Order Logic (PFHOL)

Free logic, a term coined by Lambert [19], refers to a family of logics that are free of existential presuppositions in general and with respect to the denotation of terms in particular. Terms of free logic may denote existent¹ objects, but are not necessarily required to do so. Quantification is treated as in classical logic, meaning that quantifiers range over the existing objects only. In the following, we will pursue an *inner-outer dual-domain approach* for the representation of the relationship between existing and nonexisting objects. The inner-outer dual-domain approach specifies that some domain D contains both existing and nonexisting objects, whereas the quantification domain E , a subdomain of D , contains solely the existing ones.

A free logic is known to be positive if it allows atomic formulas with terms that do not denote to be either true or false [30, 20]. For example, even though *isHuman(Pegasus)* is usually denied, *hasLegs(Pegasus)* may be regarded a valid formula since the denotation of *Pegasus* is a mythological creature that is usually depicted in the form of a winged horse (with legs).

3.1 Syntax

Except for terms, all definitions and terminology for PFHOL correspond to those presented in §2.1 for HOL. Simply typed terms of PFHOL essentially also have the same structure as terms of HOL, but we additionally include the nonlogical constant symbol $E!_{\alpha \rightarrow o} \in C_{\alpha \rightarrow o}$ in the “signature”. Apart from that, the denotation of the universal quantifier changes since quantification in free logic is traditionally limited to existing objects only. Moreover, not only quantifiers have existential import: Definite descriptions of free logic denote a unique object satisfying some property if and only if it exists and is defined [4].

3.2 Semantics

The following proposal of a positive semantics for free higher-order logic unifies two sophisticated concepts that were worked out independently by Benz Müller and Scott [8] and Farmer [12].

¹ In the paper at hand, the terms existent/existing and defined are used interchangeably even though a differentiation is advisable. The same applies to the terms nonexistent/nonexisting and undefined.

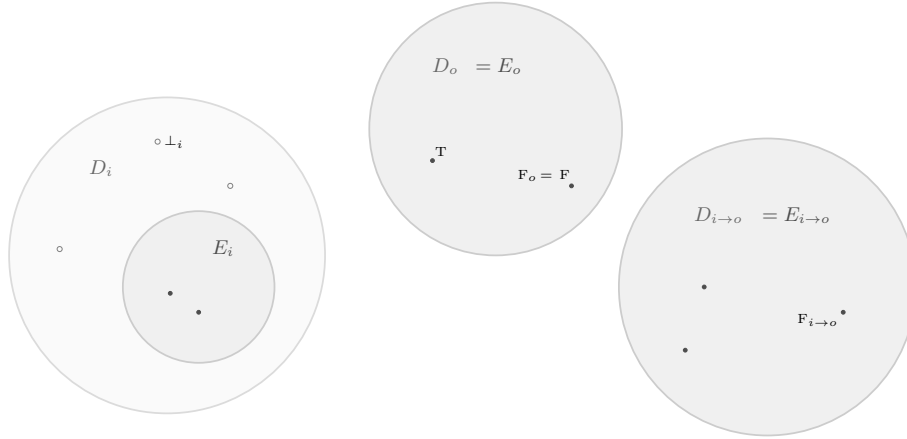


Fig. 1. Schematics of domains D_i , D_o and $D_{i \rightarrow o}$

While frames are defined exactly as in HOL, a *subframe* E is a set $\{E_\alpha : \alpha \in \mathcal{T}\}$ of nonempty sets (formally *domains*) E_α such that $E_\alpha \subset D_\alpha$ for each $\alpha \in \mathcal{T}_i$ and $E_\alpha = D_\alpha$ for each $\alpha \in \mathcal{T}_o$.² We assume, inspired by Farmer, that $\perp_\alpha \in D_\alpha \setminus E_\alpha$ for all $\alpha \in \mathcal{T}_i$ with $\perp_{\alpha \rightarrow \beta}(d) := \perp_\beta$ for all $d \in D_\alpha$. Furthermore, each set D_α with $\alpha \in \mathcal{T}_o$ contains the element F_α defined inductively by $F_o := F$ and $F_{\alpha \rightarrow \beta}(d) := F_\beta$ for all $d \in D_\alpha$. The purpose of these objects is to propagate the nondenotation of a term up through all terms containing it with \perp_i symbolizing ‘the undefinedness’ among individuals. Their intended use will be explained in the further course of this section. Exemplary schematics of the most important domains can be found in Fig. 1. A *standard model* is a triple $M = \langle D, E, I \rangle$ where D is a frame, E is a subframe and I is a family of typed interpretation functions, i.e., $I = \{I_\alpha : \alpha \in \mathcal{T}\}$. Each *interpretation function* I_α maps constants of type α to appropriate elements of D_α . The nonlogical constant $E!$ and the logical constants $=$, \neg , \vee , \forall and ι are interpreted as follows:

$$\begin{aligned}
 I(E!_{\alpha \rightarrow o}) &:= ex \in E_{\alpha \rightarrow o} && \text{s.t. for all } d \in D_\alpha, ex(d) = T \text{ iff } d \in E_\alpha \\
 I(=_{\alpha \rightarrow \alpha \rightarrow o}) &:= id \in E_{\alpha \rightarrow \alpha \rightarrow o} && \text{s.t. for all } d, d' \in D_\alpha, \\
 &&& id(d, d') = T \text{ iff } d \text{ is identical to } d' \\
 I(\neg_{o \rightarrow o}) &:= not \in E_{o \rightarrow o} && \text{s.t. } not(T) = F \text{ and } not(F) = T \\
 I(\vee_{o \rightarrow o \rightarrow o}) &:= or \in E_{o \rightarrow o \rightarrow o} && \text{s.t. } or(v_1, v_2) = T \text{ iff } v_1 = T \text{ or } v_2 = T \\
 I(\forall_{(\alpha \rightarrow o) \rightarrow o}) &:= all \in E_{(\alpha \rightarrow o) \rightarrow o} && \text{s.t. for all } f \in D_{\alpha \rightarrow o}, \\
 &&& all(f) = T \text{ iff } f(d) = T \text{ for all } d \in E_\alpha
 \end{aligned}$$

² Restricting nondenotation to the domain of individuals, i.e., to define $E_i \subset D_i$ and for all $\alpha \neq i$, $E_\alpha = D_\alpha$, is reasonable, but complicates the definition of strict functions.

$$\begin{aligned}
I(\iota_{(\alpha \rightarrow o) \rightarrow \alpha}) &:= \text{desc} \in E_{(\alpha \rightarrow o) \rightarrow \alpha} \text{ s.t. for all } f \in D_{\alpha \rightarrow o}, \\
&\quad \text{desc}(f) = d \in E_\alpha \text{ if } f(d) = \text{T} \text{ and for} \\
&\quad \text{all } d' \in E_\alpha: \text{ if } f(d') = \text{T}, \text{ then } d' = d, \\
&\quad \text{otherwise } \text{desc}(f) = \perp_\alpha \text{ if } \alpha \in \mathcal{T}_i \text{ and} \\
&\quad \text{desc}(f) = \text{F}_\alpha \text{ if } \alpha \in \mathcal{T}_o
\end{aligned}$$

The value $\llbracket s_\alpha \rrbracket^{M,g}$ of a PFHOL term s_α in a standard model M under the variable assignment g , defined in the same way as in HOL, is an object $d \in D_\alpha$ and evaluated as follows:

$$\begin{aligned}
\llbracket P_\alpha \rrbracket^{M,g} &:= I(P_\alpha) \\
\llbracket x_\alpha \rrbracket^{M,g} &:= g(x_\alpha) \\
\llbracket (s_{\alpha \rightarrow \beta} t_\alpha)_\beta \rrbracket^{M,g} &:= \llbracket s_{\alpha \rightarrow \beta} \rrbracket^{M,g} (\llbracket t_\alpha \rrbracket^{M,g}) \\
\llbracket (\lambda x_\alpha. s_\beta)_{\alpha \rightarrow \beta} \rrbracket^{M,g} &:= \text{the function } f \text{ from } D_\alpha \text{ into } D_\beta \\
&\quad \text{s.t. for all } d \in D_\alpha: f(d) = \llbracket s_\beta \rrbracket^{M,g[x \rightarrow d]}
\end{aligned}$$

The application is hereby defined in a nonstrict manner. A strict function application would be defined like this (with $\alpha \rightarrow \beta \in \mathcal{T}_i$):

$$\llbracket (s_{\alpha \rightarrow \beta} t_\alpha)_\beta \rrbracket^{M,g} := \begin{cases} \llbracket s_{\alpha \rightarrow \beta} \rrbracket^{M,g} (\llbracket t_\alpha \rrbracket^{M,g}) & \text{if } \llbracket t_\alpha \rrbracket^{M,g} \in E_\alpha^3 \\ \perp_\beta & \text{else} \end{cases}$$

A strictly applied function results in undefined if one of its arguments is undefined. In simple type theory, arguments are typically processed one after another. To be able to pass the undefined state of a once applied argument through any other possibly following arguments, the objects \perp_α were added to each relevant domain D_α . $\perp_{\alpha \rightarrow \beta}$ maps any argument of type α to \perp_β until \perp_i appears. This way, undefinedness is transmitted until the evaluation of the application has reached its end.⁴ Predicates, on the other hand, need no such special treatment. In positive free logic, (atomic) formulas denote truth or falsehood even if one of the arguments is undefined.

The definitions of truth, validity, and general validity in PFHOL are equivalent to the corresponding definitions in HOL. The partiality characteristic for free logic is implemented by a trick that exploits the objects \perp_α , which enables the functions in each domain $D_{\alpha \rightarrow \beta}$ to remain total. Hence, the generalization of standard models to Henkin models is equally applicable to PFHOL.⁵

³ Farmer also checked the function itself for existence. But since the distinction between existing and nonexisting functions – in contrast to existing/nonexisting individuals – is unusual and not well-defined, this was left out.

⁴ Restraining applications like this could lead to malformed evaluations, i.e., evaluated terms might not receive the actually intended value. For instance, the *ite* operator must be handled separately when the then- or else-parts are meant to be undefined.

⁵ As shown by Farmer and Schütte [29], it is possible to give a Henkin-style completeness proof for free higher-order logic that is defined based on a partial valuation function.

4 Embedding of PFHOL in HOL

To provide a shallow embedding of PFHOL in HOL, the “signature” of HOL has to be enriched with an additional nonlogical constant $E_{\alpha \rightarrow o} \in C_{\alpha \rightarrow o}$ denoting a unary predicate that enables an explicit distinction of existing and nonexisting objects in the domain D_α . Moreover, we assume a fixed error object e_α in each domain D_α with $\alpha \in \mathcal{T}$, which is meant to be the object that is returned by the definite description of type $(\alpha \rightarrow o) \rightarrow \alpha$ if no such object exists. We redefine the interpretation of ι thus as follows:

$$\begin{aligned} I(\iota_{(\alpha \rightarrow o) \rightarrow \alpha}) &:= \text{desc} \in D_{(\alpha \rightarrow o) \rightarrow \alpha} \text{ s.t. for all } f \in D_{\alpha \rightarrow o}, \\ &\text{desc}(f) = d \in D_\alpha \text{ if } f(d) = \text{T and for} \\ &\text{all } d' \in D_\alpha: \text{ if } f(d') = \text{T, then } d' = d, \\ &\text{otherwise } \text{desc}(f) = e_\alpha \end{aligned}$$

Obviously, for all $\alpha \in \mathcal{T}_o$: $(\forall x_\alpha. (E_{\alpha \rightarrow o} x_\alpha)_o)_o = \text{T}$, and $(E_{\alpha \rightarrow o} e_\alpha)_o = \text{F}$ for each $\alpha \in \mathcal{T}_i$. Then, a HOL term $[s_\alpha]$ is assigned to each PFHOL term s_α according to the following translation function:⁶

$$\begin{aligned} [P_\alpha] &= P_\alpha \\ [x_\alpha] &= x_\alpha \\ [(E!_{\alpha \rightarrow o} s_\alpha)_o] &= (E_{\alpha \rightarrow o} [s_\alpha])_o \\ [((=^f_{\alpha \rightarrow \alpha \rightarrow o} s_\alpha)_{\alpha \rightarrow o} t_\alpha)_o] &= ((=^h_{\alpha \rightarrow \alpha \rightarrow o} [s_\alpha])_{\alpha \rightarrow o} [t_\alpha])_o \\ [(\neg^f_{o \rightarrow o} s_o)_o] &= (\neg^h_{o \rightarrow o} [s_o])_o \\ [((\wedge^f_{o \rightarrow o \rightarrow o} s_o)_{o \rightarrow o} t_o)_o] &= ((\wedge^h_{o \rightarrow o \rightarrow o} [s_o])_{o \rightarrow o} [t_o])_o \\ [(\forall^f_{(\alpha \rightarrow o) \rightarrow o} (\lambda x_\alpha. s_o)_{\alpha \rightarrow o})_o] &= (\forall^h_{(\alpha \rightarrow o) \rightarrow o} (\lambda x_\alpha. ((E x)_o \rightarrow^h_{o \rightarrow o \rightarrow o} [s_o])_o)_{\alpha \rightarrow o})_o \\ [(\iota^f_{(\alpha \rightarrow o) \rightarrow \alpha} (\lambda x_\alpha. s_o)_{\alpha \rightarrow o})_\alpha] &= (\iota^h_{(\alpha \rightarrow o) \rightarrow \alpha} (\lambda x_\alpha. ((E x)_o \wedge^h_{o \rightarrow o \rightarrow o} [s_o])_o)_{\alpha \rightarrow o})_\alpha \\ [(s_{\alpha \rightarrow \beta} t_\alpha)_\beta] &= ([s_{\alpha \rightarrow \beta}] [t_\alpha])_\beta \\ [(\lambda x_\alpha. s_\beta)_{\alpha \rightarrow \beta}] &= (\lambda x_\alpha. [s_\beta])_{\alpha \rightarrow \beta} \end{aligned}$$

Note, that operators of HOL and PFHOL are annotated with h and f , respectively.

The main trick of this translation is that the existential import of the quantifier and description operator is secured by cleverly exploiting the additional predicate $E_{\alpha \rightarrow o}$ as a guard. When mapping definite descriptions, $[(\iota^f_{(\alpha \rightarrow o) \rightarrow \alpha} (\lambda x_\alpha. s_o)_{\alpha \rightarrow o})_\alpha]$

⁶ A similar translation, although for free first-order logic, was provided and proved to be sound and complete by Meyer and Lambert [24] and Benzmüller and Scott [8].

could also be translated into

$$\begin{aligned}
& (ite_{o \rightarrow \alpha \rightarrow \alpha \rightarrow \alpha}^{\#} \\
& \quad (\exists_{(\alpha \rightarrow o) \rightarrow o}^{\#} (\lambda x_{\alpha} \cdot \\
& \quad \quad (((E_{\alpha \rightarrow o} x_{\alpha})_o \wedge_{o \rightarrow o \rightarrow o}^{\#} [s_o])_o \\
& \quad \quad \wedge_{o \rightarrow o \rightarrow o}^{\#} (\forall_{(\alpha \rightarrow o) \rightarrow o}^{\#} (\lambda y_{\alpha} \cdot (((E_{\alpha \rightarrow o} y_{\alpha})_o \rightarrow_{o \rightarrow o \rightarrow o}^{\#} [s_o])_o \\
& \quad \quad \quad \rightarrow_{o \rightarrow o \rightarrow o}^{\#} (y_{\alpha} =_{\alpha \rightarrow \alpha \rightarrow o}^{\#} x_{\alpha})_o)_{\alpha \rightarrow o})_o)_{\alpha \rightarrow o})_o \\
& \quad \quad (\iota_{(\alpha \rightarrow o) \rightarrow \alpha}^{\#} (\lambda x_{\alpha} \cdot ((E_{\alpha \rightarrow o} x_{\alpha})_o \wedge_{o \rightarrow o \rightarrow o}^{\#} [s_o])_o)_{\alpha \rightarrow o})_{\alpha} \\
& \quad \quad e_{\alpha})_{\alpha}
\end{aligned}$$

using the if-then-else operator *ite* to ensure that the classical description definitely returns the error object e_{α} in case of no such existing object. But due to our previously done redefinition of classical definite description, this is not really necessary for this embedding. Furthermore, it is noteworthy that any term $\exists^f x. s$ is translated into $\neg^{\#} \forall^{\#} x. E x \rightarrow^{\#} \neg^{\#} s$, which is the same as $\exists^{\#} x. E x \wedge^{\#} s$.

Next, we establish the faithfulness of this embedding.

Theorem 1. $\models_{PFHOL} s_o$ if and only if $\models_{HOL} [s_o]$.

The proof of Thm. 1 is sketched in App. A; for full details see Makarenko [23].

5 Implementation in Isabelle/HOL

We have encoded the embedding from §4 in Isabelle/HOL [25]. This encoding starts out with a declaration of the base type *i* for individuals; the type *o* of HOL is associated with the predefined type `bool` in Isabelle/HOL.

```
typedecl i
```

Next, we define an existence predicate *E* for each of the base and compound types. The single quote in '*a*' indicates that this is a type variable, meaning that the definition given hereupon is polymorphic. The prefix '*f*' in this and all upcoming definitions stands for 'free'.

```
consts fExistence :: "'a  $\Rightarrow$  bool" ("E")
```

Then, we introduce a new constant *e* for every type and, in accordance with the definitions in §4, we postulate *e* of type *i* to be nonexistent and *e* of type *bool* to be `False`. Furthermore, `True` and `False` are declared as existent.

```
consts fUndef :: "'a" ("e")
axiomatization where fUndefIAxiom: " $\neg$ E (e::i)"
axiomatization where fFalsehoodBAxiom: "(e::bool) = False"
axiomatization where fTrueAxiom: "E True"
axiomatization where fFalseAxiom: "E False"
```

The embedding of the propositional logical connectives is straightforward. PFHOL operators are presented in bold-face fonts to distinguish them from HOL operators.

```
definition fIdentity :: "'a ⇒ 'a ⇒ bool" (infixr "=" 56)
  where "φ = ψ ≡ φ = ψ"
definition fNot :: "bool ⇒ bool" ("¬_" [52]53)
  where "¬φ ≡ ¬φ"
definition fOr :: "bool ⇒ bool ⇒ bool" (infixr "∨" 51)
  where "φ ∨ ψ ≡ φ ∨ ψ"
```

Now, for embedding the existential import of the universal quantifier, we utilize the existence predicate **E** of the respective type exactly as discussed in §4. Isabelle/HOL supports the introduction of syntactic sugar for binding notations, which we adopt in the following definition in order to support the more intuitive notation $\forall x. Px$ instead of writing $\forall(\lambda x. Px)$ or $\forall P$.

```
definition fForall :: "('a ⇒ bool) ⇒ bool" ("∀")
  where "∀Φ ≡ ∀x. E x → Φ x"
definition fForallBinder :: "('a ⇒ bool) ⇒ bool" (binder "∀" [8]9)
  where "∀x. φ x ≡ ∀φ"
```

For the encoding of the PFHOL operator ι , we rely on Isabelle/HOL's own definite description operator **THE**. Unlike the embedding from §4, we must here specify the exact object that will be returned if there is no unique object that has the desired property. We use Isabelle/HOL's if-then-else operator for this:

```
definition fThat :: "('a ⇒ bool) ⇒ 'a" ("I")
  where "IΦ ≡ if ∃x. E x ∧ Φ x ∧ (∀y. (E y ∧ Φ y) → (y = x))
    then THE x. E x ∧ Φ x
    else e"
definition fThatBinder :: "('a ⇒ bool) ⇒ 'a" (binder "I" [8]9)
  where "Ix. φ x ≡ Iφ"
```

We also introduced binder notation for **I**. Further PFHOL operators are embedded as follows:

```
definition fAnd :: "bool ⇒ bool ⇒ bool" (infixr "∧" 52)
  where "φ ∧ ψ ≡ ¬(¬φ ∨ ¬ψ)"
definition fImp :: "bool ⇒ bool ⇒ bool" (infixr "→" 49)
  where "φ → ψ ≡ ¬φ ∨ ψ"
definition fEquiv :: "bool ⇒ bool ⇒ bool" (infixr "↔" 50)
  where "φ ↔ ψ ≡ φ → ψ ∧ ψ → φ"
definition fExists :: "('a ⇒ bool) ⇒ bool" ("∃")
  where "∃Φ ≡ ¬(∀(λy. ¬(Φ y)))"
definition fExistsBinder :: "('a ⇒ bool) ⇒ bool" (binder "∃" [8]9)
  where "∃x. φ x ≡ ∃φ"
```

For experiments and tests, and for the Isabelle/HOL sources, see Makarenko [23], which is available at: <https://github.com/stilleben/Free-Higher-Order-Logic>.

6 Automated Assessment of Prior’s Paradox

In our practical experiments, we benefit from the fact that Isabelle/HOL integrates powerful reasoning tools such as the model finder Nitpick [10] and the meta-prover Sledgehammer [26], which invokes third-party resolution provers, SMT solvers, and higher-order provers as Satallax [11] and Leo-III [32].

Applying the meta-prover Sledgehammer together with our embedding of PFHOL in HOL onto Prior’s theorem, we end up with the following result:

```
axiomatization where fTrueAxiom: "E True"
axiomatization where fFalseAxiom: "E False"
lemma "(Q (∀p. (Q p → (¬p)))) → ((∃p. Q p ∧ p) ∧ (∃p. Q p ∧ (¬p)))"
  using Defs by (smt fFalseAxiom fTrueAxiom)
```

The theorem is valid. But as we can clearly see, the theorem is proved using the axioms `fTrueAxiom` and `fFalseAxiom` stipulating that both truth values are defined. We try it again without these:

```
lemma "(Q (∀p. (Q p → (¬p)))) → ((∃p. Q p ∧ p) ∧ (∃p. Q p ∧ (¬p)))"
  nitpick [user_axioms=true, show_all, format=2]
  oops
```

Nitpick found a counterexample for card i = 3:

```
Free variable:
  Q = (λx. _)(True := True, False := True)
Constants:
  E = (λx. _)(True := True, False := False)
  E = (λx. _)(i1 := False, i2 := False, i3 := True)
  e = i2
  e = False
```

This time the model finder Nitpick found a countermodel. Observe, that in this countermodel one of the two truth values is undefined, namely `False`. This essentially coincides with the countermodel provided by Bacon, Hawthorne, and Uzquiano. However, on a metaphysical level, it is highly questionable to shift even one of the truth values into the undefined range. Bacon et al. themselves did not find this approach for overcoming the paradox very promising and have constructed other countermodels as a substitute, which we could not reproduce with the embedding of PFHOL in HOL as presented in this paper. For these countermodels, at least three different truth values are required, and hence trivalent or other many-valued free higher-order logics should be used for that. Research has already been conducted in this direction, which, so far, has concentrated mainly on using deep embeddings [33] as opposed to adapting shallow ones [31].

An alternative option, which has been implemented and explored by Makarenko in Isabelle/HOL [23, see §5], is to embed and automate the free semantics as specifically developed by Bacon et al. to overcome this particular paradox. The semantical theory they introduce is a positive free higher-order logic based on set

theory where only (possible) worlds are taken as primitive, and the validity of propositions is then modeled as world dependent. The embedding of this ‘modal’ positive free logic has proved useful and adequate in dealing with the paradox, as was confirmed by verifying further, more reasonable countermodels to Prior’s theorem.

It is worth mentioning that there is currently a growing interest to further adapt the definitions of §3 and the embedding of §4 to develop proper notions of modal and intensional positive free higher-order logic and to embed them faithfully in HOL. An interesting application, and related ongoing work, includes the exploitation of free logic machinery in Kirchner’s embedding of hyperintensional second-order modal logic and abstract object theory in Isabelle/HOL [18, see footnote 7 and §5], which has been used for the encoding, assessment and further exploration of Zalta’s “*Principia Logico-Metaphysica*” [34].

7 Conclusion

Positive free higher-order logic and its characteristics of non-existent objects and partial functions has been faithfully represented in an adequately modified version of simple type theory. A key point of the inner-outer dual-domain approach is that partiality is only simulated instead of inherently accomodating it, such that a classical logic environment could be maintained. Subsequently, our embedding was implemented in Isabelle/HOL to support interactive and automated reasoning. We applied this embedding to Prior’s paradox and reconstructed some of the results of Bacon, Hawthorne, and Uzquiano. This shows that certain paradoxes can fruitfully be addressed in free higher-order logic. However, we were also able to verify that two-valued free logic is not enough to resolve the issue. Our ongoing research therefore also addressed other variants of free logic. Traditionally, the family of free logics involves not only positive free logic, but also negative [28], neutral [22], and supervaluational [4] free logic whose semantics differ in the way how atomic formulas with terms that do not denote are treated. Furthermore, free many-valued logic or even a free logic with more than one notion and/or degree of nonexistence could be imagined. Some of these variants have already been successfully embedded and tested in Isabelle/HOL [23], as for example negative free higher-order logic and partly also supervaluational free higher-order logic, others are still under development. Of special interest are in particular neutral free higher-order logic and, as also indicated in the previous section, many-valued (positive) free higher-order logic. Obviously, a mixture between shallow and deep embedding appears conceivable in this context and worth investigating. Fact is, nondenoting terms have always been and will always be an intriguing subject in logic, and, due to the lack of theorem provers for free logic, the development of an appropriate definition of free logic suitable for embedding in HOL as well as the automation of free logic via a semantical embedding seems more important than ever.

References

1. Peter B. Andrews. General models and extensionality. *Journal of Symbolic Logic*, 37(2):395–397, 1972.
2. Peter B. Andrews. General models, descriptions, and choice in type theory. *Journal of Symbolic Logic*, 37(2):385–394, 1972.
3. Andrew Bacon, John Hawthorne, and Gabriel Uzquiano. Higher-order free logic and the prior-kaplan paradox. *Canadian Journal of Philosophy*, 46(4-5):493–541, 2016.
4. Ermanno Bencivenga. Free logics. In Dov M. Gabbay and Franz Günthner, editors, *Handbook of Philosophical Logic. Volume III: Alternatives in Classical Logic*, pages 373–426. Springer Netherlands, Dordrecht, 1986.
5. Christoph Benzmüller. Universal (Meta-)Logical Reasoning: Recent Successes. *Science of Computer Programming*, 172:48–62, 2019.
6. Christoph Benzmüller and Peter B. Andrews. Church’s type theory. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2019 edition, 2019.
7. Christoph Benzmüller, Chad E. Brown, and Michael Kohlhase. Higher-order semantics and extensionality. *Journal of Symbolic Logic*, 69(4):1027–1088, 2004.
8. Christoph Benzmüller and Dana Scott. Automating Free Logic in HOL, with an Experimental Application in Category Theory. *Journal of Automated Reasoning*, pages 1–20, 2019.
9. Christoph Benzmüller, Xavier Parent, and Leendert van der Torre. Designing normative theories for ethical and legal reasoning: Logikey framework, methodology, and tool support. *Artificial Intelligence (accepted for publication)*, pages 1–50, 2020.
10. Jasmin Christian Blanchette and Tobias Nipkow. Nitpick: A counterexample generator for higher-order logic based on a relational model finder. In Matt Kaufmann and Lawrence C. Paulson, editors, *First International Conference on Interactive Theorem Proving*, volume 6172 of *Lecture Notes in Computer Science*, pages 131–146, Edinburgh, United Kingdom, Jul 11–14 2010. Springer Berlin, Heidelberg.
11. Chad E. Brown. Satallax: An automatic higher-order prover. In *Proceedings of the 6th International Joint Conference on Automated Reasoning, IJCAR’12*, pages 111–117, Berlin, Heidelberg, 2012. Springer.
12. William M. Farmer. A Partial Functions Version of Church’s Simple Theory of Types. *Journal of Symbolic Logic*, 55:1269–1291, 1990.
13. Solomon Feferman. Logics for termination and correctness of functional programs. In Yiannis N. Moschovakis, editor, *Logic from Computer Science*, pages 95–127, New York, NY, 1992. Springer.
14. Raymond D. Gumb. Free logic in program specification and verification. In Edgar Morscher and Alexander Hieke, editors, *New Essays in Free Logic. In Honour of Karel Lambert*, volume 23, page 157–93. Springer Netherlands, Dordrecht, 2001.
15. Raymond D. Gumb and Karel Lambert. Definitions in nonstrict positive free logic. *Modern Logic*, 7(1):25–55, 1997.
16. Leon Henkin. Completeness in the theory of types. *Journal of Symbolic Logic*, 15(2):81–91, 1950.
17. David Kaplan. A problem in possible worlds semantics. In Walter Sinnott-Armstrong, Diana Raffman, and Nicholas Asher, editors, *Modality, Morality and Belief: Essays in Honor of Ruth Barcan Marcus*, pages 41–52. Cambridge University Press, 1995.

18. Daniel Kirchner, Christoph Benzmüller, and Edward N. Zalta. Mechanizing principia logico-metaphysica in functional type theory. *Review of Symbolic Logic*, 13(1):206–218, 2020.
19. Karel Lambert. The definition of $e!$ in free logic. *Abstracts: The International Congress for Logic, Methodology and Philosophy of Science*, 1960.
20. Karel Lambert. Free Logic and the Concept of Existence. *Notre Dame Journal of Formal Logic*, 1-2(8):133–144, 1967.
21. Karel Lambert. *Philosophical Applications of Free Logic*. Oxford University Press, 1991.
22. Scott Lehmann. ‘no input, no output’ logic. In Edgar Morscher and Alexander Hieke, editors, *New Essays in Free Logic. In Honour of Karel Lambert*, volume 23, pages 147–155. Springer Netherlands, Dordrecht, 2001.
23. Irina Makarenko. Free Higher-Order Logic – Notion, Definition and Embedding in HOL. Master’s thesis, Freie Universität Berlin, 2020.
24. Robert K. Meyer and Karel Lambert. Universally Free Logic and Standard Quantification Theory. *Journal of Symbolic Logic*, 33(1):8–26, 1968.
25. Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. Isabelle/HOL — A Proof Assistant for Higher-Order Logic. <http://isabelle.in.tum.de/doc/tutorial.pdf>, 2019. Last accessed on December 30, 2019.
26. Lawrence C. Paulson and Jasmin Christian Blanchette. Three years of experience with sledgehammer, a practical link between automatic and interactive theorem provers. *Proceedings of the 8th International Workshop on the Implementation of Logics*, pages 131–146, 2015.
27. Arthur N. Prior. On a family of paradoxes. *Notre Dame Journal of Formal Logic*, 2(1):16–32, 1961.
28. Rolf Schock. *Logics Without Existence Assumptions*. Almqvist & Wiksell, Stockholm, 1968.
29. Kurt Schütte. Syntactical and semantical properties of simple type theory. *Journal of Symbolic Logic*, 25(4):305–326, 1960.
30. Dana Scott. Existence and description in formal logic. In Ralph Schoenman, editor, *Bertrand Russell: Philosopher of the Century*, pages 181–200. Little, Brown & Company, Boston, 1967. Repr. in [21], pp. 28–48.
31. Alexander Steen and Christoph Benzmüller. Sweet sixteen: Automation via embedding into classical higher-order logic. *Logic and Logical Philosophy*, 25(4):535–554, 2016.
32. Alexander Steen and Christoph Benzmüller. The higher-order prover Leo-III. In Didier Galmiche, Stephan Schulz, and Roberto Sebastiani, editors, *Automated Reasoning. IJCAR 2018*, volume 10900 of *Lecture Notes in Computer Science*, pages 108–116. Springer, 2018.
33. Jørgen Villadsen and Anders Schlichtkrull. Formalization of many-valued logics. In Henning Christiansen, M. Dolores Jiménez-López, Roussanka Loukanova, and Lawrence S. Moss, editors, *Partiality and Underspecification in Information, Languages, and Knowledge*, pages 219–256. Cambridge Scholars Press, 2017.
34. Edward N. Zalta. Principia Logico-Metaphysica. <http://mally.stanford.edu/principia.pdf>. [Draft/Excerpt; accessed: May 31, 2020].

Appendix

A Proof of Theorem 1

For the proof of the theorem, we first need to elaborate how to transform a PFHOL model M into a HOL model M^* , and a PFHOL variable assignment g into a HOL variable assignment g^* . We assume, that $D_\alpha^* = D_\alpha$ and $C_\alpha^* \setminus \{E_{\alpha \rightarrow o}\} = C_\alpha \setminus \{E!_{\alpha \rightarrow o}\}$ for all $\alpha \in \mathcal{T}$, and set $e_\alpha = \perp_\alpha$ for each $\alpha \in \mathcal{T}_i$ and $e_\alpha = F_\alpha$ for each $\alpha \in \mathcal{T}_o$. Then, $M = \langle D, E, I \rangle$ corresponds to the model $M^* = \langle D^*, I^* \rangle$ where I^* is a family of interpretation functions that assigns the standard interpretation to the logical constants $=, \neg, \vee, \forall$ and ι of HOL as described in §2. For all other constants $P_\alpha \neq E_{\alpha \rightarrow o}$, $P_\alpha \in C_\alpha^* : I^*(P_\alpha) = I(P_\alpha)$. The nonlogical constant $E_{\alpha \rightarrow o} \in C_\alpha^*$ is interpreted as follows:

$$I^*(E_{\alpha \rightarrow o}) \quad := \quad ex \quad \in D_{\alpha \rightarrow o}^* \quad \text{s.t. for all } d \in D_\alpha^*, ex(d) = T \text{ iff } d \in E_\alpha$$

We further assume $V_\alpha^* = V_\alpha$ for all $\alpha \in \mathcal{T}$, and hence, for all $x_\alpha \in V_\alpha^*$ and $\alpha \in \mathcal{T}$, $g_\alpha^*(x_\alpha) = g_\alpha(x_\alpha)$.

Next, we first need to establish the following lemma.

Lemma 1. *For all PFHOL models M and PFHOL variable assignments g ,*

$$\llbracket s_\alpha \rrbracket^{M,g} = \llbracket [s_\alpha] \rrbracket^{M^*,g^*}.$$

The detailed proof of this lemma can be found in Makarenko's thesis [23].

Theorem 1. $\vDash_{PFHOL} s_o$ *if and only if* $\vDash_{HOL} [s_o]$.

Proof.

(\rightarrow) The proof is by contraposition:

Assume $\not\vDash_{PFHOL} s_o$. Then, there exists a PFHOL model M and a variable assignment g such that $\llbracket s_o \rrbracket^{M,g} = F$. By Lemma 1, $\llbracket s_o \rrbracket^{M,g} = \llbracket [s_o] \rrbracket^{M^*,g^*} = F$. Hence, $\not\vDash_{HOL} [s_o]$.

(\leftarrow) Analogous to above by contraposition and Lemma 1.

Therefore, the embedding of PFHOL in HOL is sound and complete.