

Faithful Logic Embeddings in HOL — Deep and Shallow (Isabelle/HOL dataset)

Christoph Benzmüller

June 4, 2025

Abstract

A recipe for the simultaneous deployment of different forms of deep and shallow embeddings of non-classical logics in classical higher-order logic is presented, which enables interactive or even automated faithfulness proofs between the logic embeddings. The approach, which is particularly fruitful for logic education, is explained in detail in an associated CADE conference paper. This paper presents the corresponding Isabelle/HOL dataset (which is only slightly modified to meet AFP requirements).

Contents

1	Introduction	1
2	Preliminaries	2
3	Deep embedding of PML in HOL	3
4	Shallow embedding of PML in HOL (maximal)	3
5	Shallow embedding of PML in HOL (minimal)	4
6	Automated faithfulness proofs	5
7	Appendix: proof automation tests	6
7.1	Tests with the deep embedding	6
7.2	Tests with the maximal shallow embedding	8
7.3	Tests with the minimal shallow embedding	10

1 Introduction

The Isabelle/HOL dataset associated with [1] is presented. Sections 3, 4 and 5 present deep, maximally shallow, and minimally shallow embeddings

of propositional modal logic (PML) in classical higher-order logic (HOL). These are connected, as a novel contribution, by automated faithfulness proofs given in Sect. 6. This connection ensures that these deep and shallow embeddings can now be used interchangeably in subsequent applications. Several experiments with the presented embeddings are presented in Sect. 7. The presented work is conceptual in nature and can be adapted to other non-classical logics. For more detailed explanations of the presented material, including a discussion of related works, see [1].

2 Preliminaries

The following preliminaries are shared between all embeddings introduced in the remainder of this paper.

```
theory PMLinHOL-preliminaries
imports Main
begin

— Type declarations common for both the deep and shallow embedding
typedecl w — Type for possible worlds
typedecl S — Type for propositional constant symbols
consts p::S q::S r::S — Some propositional constant symbols
type-synonym W = w⇒bool — Type for sets of possible worlds
type-synonym R = w⇒w⇒bool — Type for accessibility relations
type-synonym V = S⇒w⇒bool — Type for valuation functions

— Some useful predicates for accessibility relations
abbreviation(input) reflexive ≡ λR::R. ∀ x. R x x
abbreviation(input) symmetric ≡ λR::R. ∀ x y. R x y → R y x
abbreviation(input) transitive ≡ λR::R. ∀ x y z. (R x y ∧ R y z) → R x z
abbreviation(input) equivrel ≡ λR::R. reflexive R ∧ symmetric R ∧ transitive R
abbreviation(input) irreflexive ≡ λR::R. ∀ x. ¬R x x
abbreviation(input) euclidean ≡ λR::R. ∀ x y z. R x y ∧ R x z → R y z
abbreviation(input) wellfounded ≡ λR::R. ∀ P::W. (∀ x. (∀ y. R y x → P y) → P x) → (∀ x. P x)
abbreviation(input) converserel ≡ λR::R. λy::w. λx::w. R x y
abbreviation(input) conversewf ≡ λR::R. wellfounded (converserel R)

— Bounded universal quantifier: ∀ x:W. φ stands for ∀ x. W x → φ x
abbreviation(input) BoundedAll::W⇒W⇒bool where BoundedAll W φ ≡ ∀ x.
W x → φ x
syntax -BoundedAll:: pttrn⇒W⇒bool⇒bool ((3∀ (-/:-)./-) [0, 0, 10] 10)
translations ∀ x:W. φ ≈ CONST BoundedAll W (λx. φ)

— Backward implication; useful for aesthetic reasons
abbreviation(input) Bimp (infixr ← 50) where φ ← ψ ≡ ψ → φ

— Some further settings
```

```

declare[[syntax-ambiguity-warning=false]]
nitpick-params[user-axioms,expect=genuine]
end

```

3 Deep embedding of PML in HOL

```

theory PMLinHOL-deep
imports PMLinHOL-preliminaries
begin

— Deep embedding (of propositional modal logic in HOL)
datatype PML = AtmD S (-d) | NotD PML (-d) | ImpD PML PML (infixr ⊃d
93) | BoxD PML (□d)

— Further logical connectives as definitions
definition OrD (infixr ∨d 92) where φ ∨d ψ ≡ ¬dφ ⊃d ψ
definition AndD (infixr ∧d 95) where φ ∧d ψ ≡ ¬d(φ ⊃d ¬dψ)
definition DiaD (◊d-) where ◊dφ ≡ ¬d(□d(¬dφ))
definition TopD (⊤d) where ⊤d ≡ pd ⊃d pd
definition BotD (⊥d) where ⊥d ≡ ¬d ⊤d

— Definition of truth of a formula relative to a model ⟨W,R,V⟩ and possible world
w
primrec RelativeTruthD :: W⇒R⇒V⇒w⇒PML⇒bool ((⟨-, -, -⟩, -) ⊨d -) where
  ⟨W,R,V⟩, w ⊨d ad = (V a w)
  | ⟨W,R,V⟩, w ⊨d ¬dφ = (¬(⟨W,R,V⟩, w ⊨d φ))
  | ⟨W,R,V⟩, w ⊨d φ ⊃d ψ = ((⟨W,R,V⟩, w ⊨d φ → ⟨W,R,V⟩, w ⊨d ψ)
  | ⟨W,R,V⟩, w ⊨d □dφ = (forall v:W. R w v → ⟨W,R,V⟩, v ⊨d φ))

— Definition of validity
definition ValD (⊨d -) where (⊨d φ) ≡ (forall W R V. ∀ w:W. ⟨W,R,V⟩, w ⊨d φ)

— Collection of definitions in a bag called DefD
named-theorems DefD declare OrD-def[DefD,simp] AndD-def[DefD,simp] DiaD-def[DefD,simp]
TopD-def[DefD,simp] BotD-def[DefD,simp] RelativeTruthD-def[DefD,simp] ValD-def[DefD,simp]
end

```

4 Shallow embedding of PML in HOL (maximal)

```

theory PMLinHOL-shallow
imports PMLinHOL-preliminaries
begin

— Shallow embedding (of propositional modal logic in HOL)
type-synonym σ = W⇒R⇒V⇒w⇒bool
definition AtmS::S⇒σ (-s) where as ≡ λ W R V w. V a w
definition NegS::σ⇒σ (-s) where ¬s φ ≡ λ W R V w. ¬(φ W R V w)
definition ImpS::σ⇒σ⇒σ (infixr ⊃s 93) where φ ⊃s ψ ≡ λ W R V w. (φ W R
V w) ⊨d ψ

```

$V w) \longrightarrow (\psi W R V w)$
definition $BoxS::\sigma \Rightarrow \sigma (\square^s)$ **where** $\square^s \varphi \equiv \lambda W R V w. \forall v:W. R w v \longrightarrow (\varphi W R V v)$

— Further logical connectives as definitions

definition OrS (**infixr** \vee^s 92) **where** $\varphi \vee^s \psi \equiv \neg^s \varphi \supset^s \psi$
definition $AndS$ (**infixr** \wedge^s 95) **where** $\varphi \wedge^s \psi \equiv \neg^s (\varphi \supset^s \neg^s \psi)$
definition $DiaS$ (\diamond^s) **where** $\diamond^s \varphi \equiv \neg^s (\square^s (\neg^s \varphi))$
definition $TopS$ (\top^s) **where** $\top^s \equiv p^s \supset^s p^s$
definition $BotS$ (\perp^s) **where** $\perp^s \equiv \neg^s \top^s$

— Definition of truth of a formula relative to a model $\langle W, R, V \rangle$ and possible world

w
definition $RelativeTruthS::\mathcal{W} \Rightarrow \mathcal{R} \Rightarrow \mathcal{V} \Rightarrow w \Rightarrow \sigma \Rightarrow \text{bool} (\langle \cdot, \cdot, \cdot \rangle, \models^s \cdot)$ **where** $\langle W, R, V \rangle, w \models^s \varphi \equiv \varphi W R V w$

— Definition of validity

definition $ValS$ ($\models^s \cdot$) **where** $\models^s \varphi \equiv \forall W R V. \forall w:W. \langle W, R, V \rangle, w \models^s \varphi$

— Collection of definitions in a bag called DefS

named-theorems $DefS$ **declare** $AtmS\text{-def}[DefS, \text{simp}]$ $NegS\text{-def}[DefS, \text{simp}]$ $ImpS\text{-def}[DefS, \text{simp}]$
 $BoxS\text{-def}[DefS, \text{simp}]$ $OrS\text{-def}[DefS, \text{simp}]$ $AndS\text{-def}[DefS, \text{simp}]$ $DiaS\text{-def}[DefS, \text{simp}]$
 $TopS\text{-def}[DefS, \text{simp}]$ $BotS\text{-def}[DefS, \text{simp}]$ $RelativeTruthS\text{-def}[DefS, \text{simp}]$ $ValS\text{-def}[DefS, \text{simp}]$
end

5 Shallow embedding of PML in HOL (minimal)

theory $PMLinHOL\text{-shallow-minimal}$
imports $PMLinHOL\text{-preliminaries}$
begin

— The accessibility relation R and the valuation function V are introduced as constants at the meta-level HOL

consts $R::\mathcal{R}$ $V::\mathcal{V}$

— Shallow embedding (of propositional modal logic in HOL)

type-synonym $\sigma = w \Rightarrow \text{bool}$
definition $AtmM::\mathcal{S} \Rightarrow \sigma (\cdot^m)$ **where** $a^m \equiv \lambda w. V a w$
definition $NegM::\sigma \Rightarrow \sigma (\neg^m)$ **where** $\neg^m \varphi \equiv \lambda w. \neg \varphi w$
definition $ImpM::\sigma \Rightarrow \sigma \Rightarrow \sigma$ (**infixr** \supset^m 93) **where** $\varphi \supset^m \psi \equiv \lambda w. \varphi w \longrightarrow \psi w$
definition $BoxM::\sigma \Rightarrow \sigma (\square^m)$ **where** $\square^m \varphi \equiv \lambda w. \forall v. R w v \longrightarrow \varphi v$

— Further logical connectives as definitions

definition OrM (**infixr** \vee^m 92) **where** $\varphi \vee^m \psi \equiv \neg^m \varphi \supset^m \psi$
definition $AndM$ (**infixr** \wedge^m 95) **where** $\varphi \wedge^m \psi \equiv \neg^m (\varphi \supset^m \neg^m \psi)$
definition $DiaM$ (\diamond^m) **where** $\diamond^m \varphi \equiv \neg^m (\square^m (\neg^m \varphi))$
definition $TopM$ (\top^m) **where** $\top^m \equiv p^m \supset^m p^m$
definition $BotM$ (\perp^m) **where** $\perp^m \equiv \neg^m \top^m$

— Definition of truth of a formula relative to a model $\langle W, R, V \rangle$ and a possible world w

definition *RelativeTruthM*:: $w \Rightarrow \sigma \Rightarrow \text{bool}$ ($\dashv^m \cdot$) **where** $w \models^m \varphi \equiv \varphi w$

— Definition of validity

definition *ValM* ($\models^m \cdot$) **where** $\models^m \varphi \equiv \forall w::w. w \models^m \varphi$

— Collection of definitions in a bag called DefM

named-theorems *DefM* **declare** *AtmM-def*[*DefM,simp*] *NegM-def*[*DefM,simp*]
ImpM-def[*DefM,simp*] *BoxM-def*[*DefM,simp*] *OrM-def*[*DefM,simp*] *AndM-def*[*DefM,simp*]
DiaM-def[*DefM,simp*] *TopM-def*[*DefM,simp*] *BotM-def*[*DefM,simp*] *RelativeTruthM-def*[*DefM,simp*]
ValM-def[*DefM,simp*]

end

6 Automated faithfulness proofs

theory *PMLinHOL-faithfulness*

imports *PMLinHOL-deep* *PMLinHOL-shallow* *PMLinHOL-shallow-minimal*
begin

— Mappings: deep to maximal shallow and deep to minimal shallow

primrec *DpToShMax* ((\dashv)) **where** $(\varphi^d) = \varphi^s \mid (\neg^d \varphi) = \neg^s (\varphi) \mid (\varphi \supset^d \psi) = (\varphi) \supset^s (\psi) \mid (\Box^d \varphi) = \Box^s (\varphi)$
primrec *DpToShMin* ([\cdot]) **where** $[\varphi^d] = \varphi^m \mid [\neg^d \varphi] = \neg^m [\varphi] \mid [\varphi \supset^d \psi] = [\varphi] \supset^m [\psi] \mid [\Box^d \varphi] = \Box^m [\varphi]$

— Proving faithfulness between deep and maximal shallow

theorem *Faithful1a*: $\forall W R V. \forall w:W. \langle W, R, V \rangle, w \models^d \varphi \longleftrightarrow \langle W, R, V \rangle, w \models^s (\varphi)$ **apply** *induct by auto*
theorem *Faithful1b*: $\models^d \varphi \longleftrightarrow \models^s (\varphi)$ **using** *Faithful1a* **by auto**

— Proving faithfulness between deep and minimal shallow

theorem *Faithful2*: $\forall w. \langle (\lambda x::w. \text{True}), R, V \rangle, w \models^d \varphi \longleftrightarrow w \models^m [\varphi]$ **apply** *induct by auto*

— Proving faithfulness maximal shallow and minimal shallow

theorem *Faithful3*: $\forall w. \langle (\lambda x::w. \text{True}), R, V \rangle, w \models^s (\varphi) \longleftrightarrow w \models^m [\varphi]$ **apply** *induct by auto*

— Additional check for soundness for the minimal shallow embedding

lemma *Sound1*: $\models^m \psi \longrightarrow (\exists \varphi. \psi = [\varphi] \wedge \models^d \varphi)$ — sledgehammer: Proof found;
metis reconstruction timeout **oops**

lemma *Sound2*: $\models^m \psi \longrightarrow (\exists \varphi. \psi = [\varphi] \wedge \models^m [\varphi])$ — sledgehammer: Proof found;
metis reconstruction timeout **oops**

end

7 Appendix: proof automation tests

7.1 Tests with the deep embedding

```
theory PMLinHOL-deep-tests
  imports PMLinHOL-deep
begin
```

— Hilbert calculus: proving that the schematic axioms and rules implied by the embedding

```
lemma H1:  $\models^d \varphi \supset^d (\psi \supset^d \varphi)$  by auto
lemma H2:  $\models^d (\varphi \supset^d (\psi \supset^d \gamma)) \supset^d ((\varphi \supset^d \psi) \supset^d (\varphi \supset^d \gamma))$  by auto
lemma H3:  $\models^d (\neg^d \varphi \supset^d \neg^d \psi) \supset^d (\psi \supset^d \varphi)$  by auto
lemma MP:  $\models^d \varphi \implies \models^d (\varphi \supset^d \psi) \implies \models^d \psi$  by auto
```

— Reasoning with the Hilbert calculus: interactive and fully automated

```
lemma HCderived1:  $\models^d (\varphi \supset^d \varphi)$  — sledgehammer(HC1 HC2 HC3 MP) returns:  
by (metis HC1 HC2 MP)
```

proof —

```
have 1:  $\models^d \varphi \supset^d ((\psi \supset^d \varphi) \supset^d \varphi)$  using H1 by auto
have 2:  $\models^d (\varphi \supset^d ((\psi \supset^d \varphi) \supset^d \varphi)) \supset^d ((\varphi \supset^d (\psi \supset^d \varphi)) \supset^d (\varphi \supset^d \varphi))$  using  
H2 by auto
have 3:  $\models^d (\varphi \supset^d (\psi \supset^d \varphi)) \supset^d (\varphi \supset^d \varphi)$  using 1 2 MP by meson
have 4:  $\models^d \varphi \supset^d (\psi \supset^d \varphi)$  using H1 by auto
thus ?thesis using 3 4 MP by meson
qed
```

```
lemma HCderived2:  $\models^d \varphi \supset^d (\neg^d \varphi \supset^d \psi)$  by (metis H1 H2 H3 MP)
lemma HCderived3:  $\models^d (\neg^d \varphi \supset^d \varphi) \supset^d \varphi$  by (metis H1 H2 H3 MP)
lemma HCderived4:  $\models^d (\varphi \supset^d \psi) \supset^d (\neg^d \psi \supset^d \neg^d \varphi)$  by auto
```

— Modal logic: the schematic necessitation rule and distribution axiom are implied

```
lemma Nec:  $\models^d \varphi \implies \models^d \Box^d \varphi$  by auto
lemma Dist:  $\models^d \Box^d (\varphi \supset^d \psi) \supset^d (\Box^d \varphi \supset^d \Box^d \psi)$  by auto
```

— Correspondence theory: correct statements

```
lemma cM: reflexive R  $\longleftrightarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \Box^d \varphi \supset^d \varphi$ ) —  
sledgehammer: Proof found oops
```

```
lemma cBa: symmetric R  $\longrightarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \varphi \supset^d \Box^d (\Diamond^d \varphi)$ )  
by auto
```

```
lemma cBb: symmetric R  $\longleftarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \varphi \supset^d \Box^d (\Diamond^d \varphi)$ )  
— sledgehammer: No proof oops
```

```
lemma c4a: transitive R  $\longrightarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \Box^d \varphi \supset^d \Box^d (\Box^d \varphi)$ )  
by (metis RelativeTruthD.simps)
```

```
lemma c4b: transitive R  $\longleftarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \Box^d \varphi \supset^d \Box^d (\Box^d \varphi)$ )  
— sledgehammer: No proof oops
```

— Correspondence theory: incorrect statements

```
lemma reflexive R  $\longrightarrow$  ( $\forall \varphi W V. \forall w:W. \langle W, R, V \rangle, w \models^d \Box^d \varphi \supset^d \Box^d (\Box^d \varphi)$ )  
nitpick[card w=3] oops — nitpick: Cex.
```

— Simple, incorrect validity statements

lemma $\models^d \varphi \supset^d \square^d \varphi$ **nitpick**[*card w=2, card S= 1*] **oops** — nitpick: Counterexample: modal collapse not implied

lemma $\models^d \square^d (\square^d \varphi \supset^d \square^d \psi) \vee^d \square^d (\square^d \psi \supset^d \square^d \varphi)$ **nitpick**[*card w=3*] **oops** — nitpick: Counterexample

lemma $\models^d (\Diamond^d(\Box^d \varphi)) \supset^d \Box^d(\Diamond^d \varphi)$ **nitpick**[*card w=2*] **oops** — nitpick: Counterexample

— Implied axiom schemata in S5

lemma *KB*: *symmetric R* $\longrightarrow (\forall \varphi \psi W V. \forall w:W. \langle W,R,V \rangle, w \models^d (\Diamond^d(\Box^d \varphi)) \supset^d \Box^d(\Diamond^d \varphi))$ **by auto**

lemma *K4B*: *symmetric R* \wedge *transitive R* $\longrightarrow (\forall \varphi \psi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \Box^d(\Box^d \varphi \supset^d \Box^d \psi) \vee^d \Box^d(\Box^d \psi \supset^d \Box^d \varphi))$ **by (smt OrD-def RelativeTruthD.simps)**

end

theory *PMLinHOL-deep-further-tests*
imports *PMLinHOL-deep-tests*
begin

— Implied modal principle

lemma *K-Dia*: $\models^d (\Box^d(\varphi \supset^d \psi)) \supset^d ((\Diamond^d \varphi) \supset^d \Diamond^d \psi)$ **by auto**

— Example 6.10 of Sider (2009) Logic for Philosophy

lemma *T1a*: $\models^d \Box^d p^d \supset^d ((\Diamond^d q^d) \supset^d \Diamond^d(p^d \wedge^d q^d))$ **by auto** — fast automation in meta-logic HOL

lemma *T1b*: $\models^d \Box^d p^d \supset^d ((\Diamond^d q^d) \supset^d \Diamond^d(p^d \wedge^d q^d))$ — alternative interactive proof in modal object logic K

proof —

have 1: $\models^d p^d \supset^d (q^d \supset^d (p^d \wedge^d q^d))$ **unfolding AndD-def using H1 H2 H3 MP by metis**

have 2: $\models^d \Box^d(p^d \supset^d (q^d \supset^d (p^d \wedge^d q^d)))$ **using 1 Nec by metis**

have 3: $\models^d \Box^d p^d \supset^d (q^d \supset^d (p^d \wedge^d q^d))$ **using 2 Dist MP by metis**

have 4: $\models^d (\Box^d(q^d \supset^d (p^d \wedge^d q^d))) \supset^d ((\Diamond^d q^d) \supset^d \Diamond^d(p^d \wedge^d q^d))$ **using K-Dia by metis**

have 5: $\models^d \Box^d p^d \supset^d ((\Diamond^d q^d) \supset^d \Diamond^d(p^d \wedge^d q^d))$ **using 3 4 H1 H2 MP by metis**

thus ?thesis .

qed

end

theory *PMLinHOL-deep-Loeb-tests*
imports *PMLinHOL-deep*
begin

— Löb axiom: with the deep embedding automated reasoning tools are not very responsive

lemma *Loeb1*: $\forall \varphi. \models^d \Box^d(\Box^d \varphi \supset^d \varphi) \supset^d \Box^d \varphi$ **nitpick**[*card w=1, card S=1*] **oops** — nitpick: Counterexample

lemma *Loeb2*: (*conversewf R* \wedge *transitive R*) \longrightarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \square^d(\square^d\varphi \supset^d \varphi) \supset^d \square^d(\square^d\varphi)$) — sledgehammer: No Proof **oops**
lemma *Loeb3*: (*conversewf R* \wedge *transitive R*) \longleftarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \square^d(\square^d\varphi \supset^d \varphi) \supset^d \square^d(\square^d\varphi)$) — sledgehammer: No Proof **oops**
lemma *Loeb3a*: *conversewf R* \longleftarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \square^d(\square^d\varphi \supset^d \varphi) \supset^d \square^d(\square^d\varphi)$) — sledgehammer: No Proof **oops**
lemma *Loeb3b*: *transitive R* \longleftarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \square^d(\square^d\varphi \supset^d \varphi) \supset^d \square^d(\square^d\varphi)$) — sledgehammer: No Proof **oops**
lemma *Loeb3c*: *irreflexive R* \longleftarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^d \square^d(\square^d\varphi \supset^d \varphi) \supset^d \square^d(\square^d\varphi)$) — sledgehammer: No Proof **oops**
end

7.2 Tests with the maximal shallow embedding

theory *PMLinHOL-shallow-tests*

imports *PMLinHOL-shallow*

begin

— Hilbert calculus: proving that the schematic axioms and rules implied by the embedding

lemma *H1*: $\models^s \varphi \supset^s (\psi \supset^s \varphi)$ **by auto**
lemma *H2*: $\models^s (\varphi \supset^s (\psi \supset^s \gamma)) \supset^s ((\varphi \supset^s \psi) \supset^s (\varphi \supset^s \gamma))$ **by auto**
lemma *H3*: $\models^s (\neg^s \varphi \supset^s \neg^s \psi) \supset^s (\psi \supset^s \varphi)$ **by auto**
lemma *MP*: $\models^s \varphi \implies \models^s (\varphi \supset^s \psi) \implies \models^s \psi$ **by auto**

— Reasoning with the Hilbert calculus: interactive and fully automated

lemma *HCderived1*: $\models^s (\varphi \supset^s \psi)$ — sledgehammer(HC1 HC2 HC3 MP) returns:
by (metis HC1 HC2 MP)

proof —

have 1: $\models^s \varphi \supset^s ((\psi \supset^s \varphi) \supset^s \varphi)$ **using H1 by auto**

have 2: $\models^s (\varphi \supset^s ((\psi \supset^s \varphi) \supset^s \varphi)) \supset^s ((\varphi \supset^s (\psi \supset^s \varphi)) \supset^s (\varphi \supset^s \varphi))$ **using H2 by auto**

have 3: $\models^s (\varphi \supset^s (\psi \supset^s \varphi)) \supset^s (\varphi \supset^s \varphi)$ **using 1 2 MP by meson**

have 4: $\models^s \varphi \supset^s (\psi \supset^s \varphi)$ **using H1 by auto**

thus ?*thesis* **using** 3 4 MP **by meson**

qed

lemma *HCderived2*: $\models^s \varphi \supset^s (\neg^s \varphi \supset^s \psi)$ **by (metis H1 H2 H3 MP)**

lemma *HCderived3*: $\models^s (\neg^s \varphi \supset^s \varphi) \supset^s \varphi$ **by (metis H1 H2 H3 MP)**

lemma *HCderived4*: $\models^s (\varphi \supset^s \psi) \supset^s (\neg^s \psi \supset^s \neg^s \varphi)$ **by auto**

— Modal logic: the schematic necessitation rule and distribution axiom are implied

lemma *Nec*: $\models^s \varphi \implies \models^s \square^s \varphi$ **by auto**

lemma *Dist*: $\models^s \square^s (\varphi \supset^s \psi) \supset^s (\square^s \varphi \supset^s \square^s \psi)$ **by auto**

— Correspondence theory: correct statements

lemma *cM:reflexive R* \longleftrightarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s \varphi \supset^s \varphi$) — sledgehammer: Proof found **oops**

lemma *cBa: symmetric R* \longrightarrow ($\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \varphi \supset^s \square^s (\diamond^s \varphi)$)

by auto

lemma *cBb: symmetric R* $\leftarrow (\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \varphi \supset^s \square^s(\diamond^s \varphi))$
— sledgehammer: No proof **oops**

lemma *c4a: transitive R* $\longrightarrow (\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s \varphi \supset^s \square^s(\square^s \varphi))$
by (smt DefS)

lemma *c4b: transitive R* $\leftarrow (\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s \varphi \supset^s \square^s(\square^s \varphi))$
— sledgehammer: No proof **oops**

— Correspondence theory: incorrect statements

lemma *reflexive R* $\longrightarrow (\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s \varphi \supset^s \square^s(\square^s \varphi))$
nitpick[*card w=3*] **oops** — nitpick: Counterexample

— Simple, incorrect validity statements

lemma $\models^s \varphi \supset^s \square^s \varphi$ **nitpick**[*card w=2, card S= 1*] **oops** — nitpick: Counterexample: modal collapse not implied

lemma $\models^s \square^s(\square^s \varphi \supset^s \square^s \psi) \vee^s \square^s(\square^s \psi \supset^s \square^s \varphi)$ **oops** — nitpick[*card w=3*]
returns: unknown

lemma $\models^s (\diamond^s(\square^s \varphi)) \supset^s \square^s(\diamond^s \varphi)$ **nitpick**[*card w=2*] **oops** — nitpick: Counterexample

— Implied axiom schemata in S5

lemma *KB: symmetric R* $\longrightarrow (\forall \varphi \psi W V. \forall w:W. \langle W,R,V \rangle, w \models^s (\diamond^s(\square^s \varphi)) \supset^s \square^s(\diamond^s \varphi))$ **by auto**

lemma *K4B: symmetric R \wedge transitive R* $\longrightarrow (\forall \varphi \psi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s \varphi \supset^s \square^s \psi) \vee^s \square^s(\square^s \psi \supset^s \square^s \varphi))$ **by** (smt DefS)
end

theory *PMLinHOL-shallow-further-tests*

imports *PMLinHOL-shallow-tests*
begin

— Implied modal principle

lemma *K-Dia: $\models^s (\square^s(\varphi \supset^s \psi)) \supset^s ((\diamond^s \varphi) \supset^s \diamond^s \psi)$* **by auto**

— Example 6.10 of Sider (2009) Logic for Philosophy

lemma *T1a: $\models^s \square^s p^s \supset^s ((\diamond^s q^s) \supset^s (p^s \wedge^s q^s))$* **by auto** — fast automation in meta-logic HOL

lemma *T1b: $\models^s \square^s p^s \supset^s ((\diamond^s q^s) \supset^s \diamond^s (p^s \wedge^s q^s))$* — alternative interactive proof in modal object logic K

proof —

have 1: $\models^s p^s \supset^s (q^s \supset^s (p^s \wedge^s q^s))$ **unfolding** AndS-def **using** H1 H2 H3
MP by metis

have 2: $\models^s \square^s(p^s \supset^s (q^s \supset^s (p^s \wedge^s q^s)))$ **using** 1 Nec **by metis**

have 3: $\models^s \square^s p^s \supset^s \square^s(q^s \supset^s (p^s \wedge^s q^s))$ **using** 2 Dist MP **by metis**

have 4: $\models^s (\square^s(q^s \supset^s (p^s \wedge^s q^s))) \supset^s ((\diamond^s q^s) \supset^s \diamond^s(p^s \wedge^s q^s))$ **using** K-Dia

by metis

have 5: $\models^s \square^s p^s \supset^s ((\diamond^s q^s) \supset^s \diamond^s (p^s \wedge^s q^s))$ **using** 3 4 H1 H2 MP **by metis**

thus ?thesis .

```
qed
end
```

```
theory PMLinHOL-shallow-Loeb-tests
imports PMLinHOL-shallow
begin
```

— Löb axiom: with the minimal shallow embedding automated reasoning tools are still partly responsive

```
lemma Loeb1:  $\forall \varphi. \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$  nitpick[card w=1,card S=1] oops
— nitpick: Counterexample
```

```
lemma Loeb2: (conversewf R ∧ transitive R) —> ( $\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$ ) — sledgehammer: Proof found oops
```

```
lemma Loeb3: (conversewf R ∧ transitive R) —> ( $\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$ ) — sledgehammer: No Proof oops
```

```
lemma Loeb3a: conversewf R —> ( $\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$ ) — sledgehammer: Proof found oops
```

```
lemma Loeb3b: transitive R —> ( $\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$ ) — sledgehammer: No Proof oops
```

```
lemma Loeb3c: irreflexive R —> ( $\forall \varphi W V. \forall w:W. \langle W,R,V \rangle, w \models^s \square^s(\square^s\varphi \supset^s \varphi) \supset^s \square^s\varphi$ ) — sledgehammer: Proof found oops
```

```
end
```

7.3 Tests with the minimal shallow embedding

```
theory PMLinHOL-shallow-minimal-tests
imports PMLinHOL-shallow-minimal
begin
```

— Hilbert calculus: proving that the schematic axioms and rules implied by the embedding

```
lemma H1:  $\models^m \varphi \supset^m (\psi \supset^m \varphi)$  by auto
lemma H2:  $\models^m (\varphi \supset^m (\psi \supset^m \gamma)) \supset^m ((\varphi \supset^m \psi) \supset^m (\varphi \supset^m \gamma))$  by auto
lemma H3:  $\models^m (\neg^m \varphi \supset^m \neg^m \psi) \supset^m (\psi \supset^m \varphi)$  by auto
lemma MP:  $\models^m \varphi \Rightarrow \models^m (\varphi \supset^m \psi) \Rightarrow \models^m \psi$  by auto
```

— Reasoning with the Hilbert calculus: interactive and fully automated

```
lemma HCderived1:  $\models^m (\varphi \supset^m \varphi)$  — sledgehammer(HC1 HC2 HC3 MP) returns:  
by (metis HC1 HC2 MP)
```

proof —

```
have 1:  $\models^m \varphi \supset^m ((\psi \supset^m \varphi) \supset^m \varphi)$  using H1 by auto
```

```
have 2:  $\models^m (\varphi \supset^m ((\psi \supset^m \varphi) \supset^m \varphi)) \supset^m ((\varphi \supset^m (\psi \supset^m \varphi)) \supset^m (\varphi \supset^m \varphi))$ 
```

using H2 by auto

```
have 3:  $\models^m (\varphi \supset^m (\psi \supset^m \varphi)) \supset^m (\varphi \supset^m \varphi)$  using 1 2 MP by meson
```

```
have 4:  $\models^m \varphi \supset^m (\psi \supset^m \varphi)$  using H1 by auto
```

```
thus ?thesis using 3 4 MP by meson
```

qed

```
lemma HCderived2:  $\models^m \varphi \supset^m (\neg^m \varphi \supset^m \psi)$  by (metis H1 H2 H3 MP)
```

lemma *HCderived3*: $\models^m (\neg^m \varphi \supset^m \varphi) \supset^m \varphi$ **by** (*metis H1 H2 H3 MP*)
lemma *HCderived4*: $\models^m (\varphi \supset^m \psi) \supset^m (\neg^m \psi \supset^m \neg^m \varphi)$ **by** *auto*

— Modal logic: the schematic necessitation rule and distribution axiom are implied

lemma *Nec*: $\models^m \varphi \implies \models^m \Box^m \varphi$ **by** (*smt DefM*)

lemma *Dist*: $\models^m \Box^m(\varphi \supset^m \psi) \supset^m (\Box^m \varphi \supset^m \Box^m \psi)$ **by** *auto*

— Correspondence theory: correct statements

lemma *cM:reflexive* $R \longleftrightarrow (\forall \varphi. \models^m \Box^m \varphi \supset^m \varphi)$ **by** *auto*

lemma *cBa: symmetric* $R \longrightarrow (\forall \varphi. \models^m \varphi \supset^m \Box^m \Diamond^m \varphi)$ **by** *auto*

lemma *cBb: symmetric* $R \longleftarrow (\forall \varphi. \models^m \varphi \supset^m \Box^m \Diamond^m \varphi)$ **by** (*metis DefM*)

lemma *c4a: transitive* $R \longrightarrow (\forall \varphi. \models^m \Box^m \varphi \supset^m \Box^m(\Box^m \varphi))$ **by** (*smt DefM*)

lemma *c4b: transitive* $R \longleftarrow (\forall \varphi. \models^m \Box^m \varphi \supset^m \Box^m(\Box^m \varphi))$ **by** *auto*

— Correspondence theory: incorrect statements

lemma *reflexive* $R \longrightarrow (\forall \varphi. \models^m \Box^m \varphi \supset^m \Box^m(\Box^m \varphi))$ **nitpick**[*card w=3, show-all*]

oops — nitpick: Counterexample

— Simple, incorrect validity statements

lemma $\models^m \varphi \supset^m \Box^m \varphi$ **nitpick**[*card w=2, card S= 1*] **oops** — nitpick: Counterexample: modal collapse not implied

lemma $\models^m \Box^m(\Box^m \varphi \supset^m \Box^m \psi) \vee^m \Box^m(\Box^m \psi \supset^m \Box^m \varphi)$ **nitpick**[*card w=3*] **oops** — nitpick: Counterexample

lemma $\models^m (\Diamond^m(\Box^m \varphi)) \supset^m \Box^m(\Diamond^m \varphi)$ **nitpick**[*card w=2*] **oops** — nitpick: Counterexample

— Implied axiom schemata in S5

lemma *KB: symmetric* $R \longrightarrow (\forall \varphi \psi. \models^m (\Diamond^m(\Box^m \varphi)) \supset^m \Box^m(\Diamond^m \varphi))$ **by** *auto*

lemma *K4B: symmetric* $R \wedge \text{transitive } R \longrightarrow (\forall \varphi \psi. \models^m \Box^m(\Box^m \varphi \supset^m \Box^m \psi) \vee^m \Box^m(\Box^m \psi \supset^m \Box^m \varphi))$ **by** (*smt DefM*)

end

theory *PMLinHOL-shallow-minimal-further-tests*

imports *PMLinHOL-shallow-minimal-tests*

begin

— Implied modal principle

lemma *K-Dia*: $\models^m (\Box^m(\varphi \supset^m \psi)) \supset^m ((\Diamond^m \varphi) \supset^m \Diamond^m \psi)$ **by** *auto*

— Example 6.10 of Sider (2009) Logic for Philosophy

lemma *T1a*: $\models^m \Box^m p^m \supset^m ((\Diamond^m q^m) \supset^m \Diamond^m(p^m \wedge^m q^m))$ **by** *auto* — fast automation in meta-logic HOL

lemma *T1b*: $\models^m \Box^m p^m \supset^m ((\Diamond^m q^m) \supset^m \Diamond^m(p^m \wedge^m q^m))$ — alternative interactive proof in modal object logic K

proof —

have 1: $\models^m p^m \supset^m (q^m \supset^m (p^m \wedge^m q^m))$ **unfolding** *AndM-def* **using** *H1 H2 H3 MP* **by** *metis*

have 2: $\models^m \Box^m(p^m \supset^m (q^m \supset^m (p^m \wedge^m q^m)))$ **using** 1 *Nec* **by** *metis*

have 3: $\models^m \Box^m p^m \supset^m \Box^m(q^m \supset^m (p^m \wedge^m q^m))$ **using** 2 *Dist MP* **by** *metis*

```

have 4:  $\models^m (\square^m(q^m \supset^m (p^m \wedge^m q^m))) \supset^m ((\Diamond^m q^m) \supset^m \Diamond^m(p^m \wedge^m q^m))$ 
using K-Dia by metis
have 5:  $\models^m \square^m p^m \supset^m ((\Diamond^m q^m) \supset^m \Diamond^m(p^m \wedge^m q^m))$  using 3 4 H1 H2 MP
by metis
thus ?thesis .
qed
end

theory PMLinHOL-shallow-minimal-Loeb-tests
imports PMLinHOL-shallow-minimal
begin

— Löb axiom: with the minimal shallow embedding automated reasoning tools are
still partly responsive
lemma Loeb1:  $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$  nitpick[card w=1,card S=1]
oops — nitpick: Counterexample.
lemma Loeb2: (conversewf R  $\wedge$  transitive R) — $\rightarrow$  ( $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$ ) — sh: Proof found oops
lemma Loeb3: (conversewf R  $\wedge$  transitive R) — $\leftarrow$  ( $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$ ) — sh: No Proof oops
lemma Loeb3a: conversewf R — $\leftarrow$  ( $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$ ) unfolding
DefM by blast
lemma Loeb3b: transitive R — $\leftarrow$  ( $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$ ) — sledge-
hammer: No Proof oops
lemma Loeb3c: irreflexive R — $\leftarrow$  ( $\forall \varphi. \models^m \square^m(\square^m \varphi \supset^m \varphi) \supset^m \square^m \varphi$ ) — sledge-
hammer: Proof found oops
end

```

References

- [1] C. Benzmüller. Faithful logic embeddings in HOL — deep and shallow. In *Automated Deduction – CADE-30 – 30th International Conference on Automated Deduction, Stuttgart, Germany, July 28–31, Proceedings*, Lecture Notes in Computer Science. Springer, 2025. To appear (preprint: [arXiv:2502.19311v1](https://arxiv.org/abs/2502.19311v1)).