
**Explicit Fixed Points
in interpretability logic**

Dick de Jongh
&
Albert Visser

Department of Philosophy
University of Utrecht

Logic Group

Preprint Series

No. 44



Department of Philosophy

University of Utrecht

**EXPLICIT FIXED POINTS
IN INTERPRETABILITY LOGIC**

Dick de Jongh

&

Albert Visser

MARCH 1989

DEPARTMENT OF PHILOSOPHY
UNIVERSITY OF UTRECHT
Heidelberglaan 2
3584 CS UTRECHT
The Netherlands.

1 Introduction

The basic theorems of *Provability Logic* are three in number. First is the Arithmetical Completeness Theorem. The second place is shared by the theorems affirming the Uniqueness of Fixed Points and the Explicit Definability of Fixed Points. In this paper we consider the problem of Uniqueness and Explicit Definability of Fixed Points for *Interpretability Logic*. It turns out that Uniqueness is an immediate corollary of a theorem of Smoryński, so most of the paper is devoted to proving Explicit Definability. More sketchy proofs of this Explicit Definability Theorem were given in Visser[88P] and, model-theoretically, in De Jongh & Veltman[88].

Interpretability Logic results from Provability Logic by adding a Binary Modal Operator \triangleright . If T is a given theory containing enough Arithmetic, we can interpret the modal language into the language of T in the usual way. We interpret $A \triangleright B$ as: (the formalization of) $T+B$ is relatively interpretable in $T+A$. Interpretations of a modal language of this kind were first considered in Hájek[81] and Švejdar[83]. For a more extensive introduction to the various systems of Interpretability Logic see Visser[88].

The system **IL**, the basic system of Interpretability Logic considered in this paper, is a system of arithmetically valid principles. **IL** is definitely arithmetically incomplete, but very natural from the modal point of view. The language of **IL** is the usual language of Modal Propositional Logic with an extra binary connective \triangleright . The theory **IL** is given as Propositional Logic plus:

- L1 $\vdash A \Rightarrow \vdash \Box A$
- L2 $\vdash \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$
- L3 $\vdash \Box A \rightarrow \Box \Box A$
- L4 $\vdash \Box(\Box A \rightarrow A) \rightarrow \Box A$
- J1 $\vdash \Box(A \rightarrow B) \rightarrow A \triangleright B$
- J2 $\vdash (A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C$
- J3 $\vdash (A \triangleright C) \wedge (B \triangleright C) \rightarrow A \vee B \triangleright C$
- J4 $\vdash A \triangleright B \rightarrow (\Diamond A \rightarrow \Diamond B)$
- J5 $\vdash \Diamond A \triangleright A$

In the conventions for leaving out parentheses \triangleright binds stronger than \rightarrow , but less strong than the other connectives. The principle J5 is the Interpretation Existence Lemma: it is a syntactic form of the Model Existence Lemma.

L3 is doubly superfluous: as is well-known it can be derived from L4, but in **IL** it can also be derived from J4 and J5. (Interestingly, on the arithmetical side the alternative proof leads in some cases to better estimates on the length of proofs of provability.)

\mathbf{IL} is valid for arithmetical interpretations in *adequate* theories T , i.e. theories into which $\mathbf{I}\Delta_0 + \Omega_1$ is translatable and whose axiom sets can be represented by a Δ_1^b -formula (see Buss[85] for a definition of the bounded hierarchy). It is surely arithmetically incomplete: the principle W introduced immediately below and some other principles discussed in section 4 are not provable in \mathbf{IL} , but valid in every adequate theory.

Kripke models for \mathbf{IL} were invented by Frank Veltman and a Kripke model completeness theorem was proved by De Jongh & Veltman (see De Jongh & Veltman[88]).

Other important interpretability logics which have been studied are the extensions \mathbf{ILW} , \mathbf{ILP} and \mathbf{ILM} of \mathbf{IL} obtained by adding to \mathbf{IL} respectively the principles W , P , M :

$$\begin{array}{ll} W & \vdash A \triangleright B \rightarrow A \triangleright B \wedge \Box \neg A \\ P & \vdash A \triangleright B \rightarrow \Box(A \triangleright B) \\ M & \vdash A \triangleright B \rightarrow A \wedge \Box C \triangleright B \wedge \Box C \end{array}$$

Kripke model completeness theorems for \mathbf{IL} , \mathbf{ILP} and \mathbf{ILM} were proved by De Jongh & Veltman ([88]), arithmetic completeness was proved for \mathbf{ILP} by Visser ([88]) with respect to all sequential finitely axiomatizable theories extending $\mathbf{I}\Delta_0 + \text{SUPEREXP}$, and for \mathbf{ILM} arithmetic completeness with respect to \mathbf{PA} and other essentially reflexive theories has been established independently by Berarducci and Shavrukov. \mathbf{ILW} , which is contained in both \mathbf{ILP} and \mathbf{ILM} , is still arithmetically valid in any adequate theory T . It is conjectured that \mathbf{ILW} contains precisely the principles valid in every reasonable theory T , i.e.:

$$\mathbf{ILW} \vdash A \Leftrightarrow \text{for all adequate } T, \text{ for all interpretations } * \text{ in } T, T \vdash (A)^*.$$

The restriction to \mathbf{IL} is for our purpose in this paper no limitation: theories that are arithmetically complete are evidently extensions of \mathbf{IL} and every extension of \mathbf{IL} inherits Uniqueness and Explicit Definability of Fixed Points from \mathbf{IL} . In one respect restriction to \mathbf{IL} does make a difference however: in a stronger theory the explicit fixed points could take a simpler form. We show that this indeed happens for \mathbf{ILW} .

Although the Explicit Definability of Fixed Points is a beautiful property for a system to have, the other side of the coin is that fixed points of formulas expressible in a system satisfying it can never give anything new. Thus, one cannot expect in pure interpretability logic interesting fixed points like the Rosser fixed points featuring in provability logic extended with witness comparison symbols.

2 Unique & Explicit Fixed Points in general

For our purposes we need the careful discussion of bi-modal self-reference in Smoryński[85] (p.172-176) in a slightly adapted form. Let \mathbf{SR}_0 be the following system in the the language of modal propositional logic extended with a binary operator #:

- L1 $\vdash A \Rightarrow \vdash \Box A$
 L2 $\vdash \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$
 L3 $\vdash \Box A \rightarrow \Box \Box A$
 L4 $\vdash \Box(\Box A \rightarrow A) \rightarrow \Box A$
 E $\vdash \Box(A \leftrightarrow B) \rightarrow (A \# C \leftrightarrow B \# C)$
 $\vdash \Box(A \leftrightarrow B) \rightarrow (C \# A \leftrightarrow C \# B)$

Here E stands for Extensionality.

Define $\Box^+ A := (A \wedge \Box A)$. We write A_p for a formula A in which p possibly occurs, in which case, e.g., AB stands for the result of the substitution of B for p in A_p and AAB for the result of substituting AB for p in A_p . We say that p occurs *modalized* in A_p , if p occurs in A_p only in the scope of \Box and #. Two immediate consequences of our theory are the Substitution Principles S_1, S_2, S_3 and Löb's Rule LR:

- S_1 $\vdash B \leftrightarrow C \Rightarrow \vdash AB \leftrightarrow AC$
 S_2 $\vdash \Box^+(B \leftrightarrow C) \rightarrow (AB \leftrightarrow AC)$
 S_3 Suppose p is modalized in A_p , then:
 $\vdash \Box(B \leftrightarrow C) \rightarrow (AB \leftrightarrow AC)$
 LR Let B be a conjunction of formulas of the form $\Box C$ or $\Box^+ C$, then:
 $\vdash B \rightarrow (\Box A \rightarrow A) \Rightarrow \vdash B \rightarrow A$

2.1 Uniqueness Theorem

Suppose p occurs modalized in A , then: $\mathbf{SR}_0 \vdash (\Box^+(p \leftrightarrow A_p) \wedge \Box^+(q \leftrightarrow A_q)) \rightarrow (p \leftrightarrow q)$.

Proof: By S_3 : $\vdash (\Box^+(p \leftrightarrow A_p) \wedge \Box^+(q \leftrightarrow A_q)) \rightarrow (\Box(p \leftrightarrow q) \rightarrow (p \leftrightarrow q))$. So LR gives us the desired conclusion. \square

The Uniqueness Theorem was in its original form due to Bernardi, De Jongh and Sambin. In its present form it is due to Smoryński. Assuming the modal completeness theorem an alternative model-theoretic proof along the lines of the implicit definability theorem (see theorem 3.1, p.109, Smoryński[85]) is easily given.

Let SR_1 be SR_0 plus the following axiom:

$$L3' \quad \vdash A\#B \rightarrow \Box(A\#B).$$

An immediate consequence of SR_1 is LR^+ :

$$LR^+ \quad \text{Let } B \text{ be a conjunction of formulas of the form } \Box C \text{ or } \Box^+ C \text{ or } C\#D, \text{ then:}$$

$$\vdash B \rightarrow (\Box A \rightarrow A) \Rightarrow \vdash B \rightarrow A$$

In this general setting the Explicit Definability Theorem is split up into two parts, from which the theorem itself can then be deduced as a Corollary.

2.2 Explicit Definability Theorem, part 1

Let A_p be either of the form $\Box B_p$ or $B_p\#C_p$, then there is a formula D such that: $SR_1 \vdash D \leftrightarrow AD$.

Proof: Suppose A_p is $\Box B_p$ or $B_p\#C_p$. Take $D := A \top$. We have from $L3'$: $\vdash A \top \rightarrow \Box^+(A \top \leftrightarrow \top)$, and hence by S_2 : $\vdash A \top \rightarrow A A \top$. On the other hand by S_3 : $\vdash A A \top \rightarrow (\Box A \top \rightarrow A \top)$. So LR^+ gives us: $\vdash A A \top \rightarrow A \top$. \square

To state the second part of the Explicit Definability Theorem we introduce a simple notion. Fix for the moment a propositional variable p . We write:

$A_p \leq B_p$: \Leftrightarrow whenever A_p can be written as $A^*(p, E_1 q, \dots, E_n q)$, where q does not occur in $A^*(p, r_1, \dots, r_n)$ and p does not occur in the $E_k q$, then B_p can be written as $B^*(p, E_1 q, \dots, E_n q)$, where q does not occur in $B^*(p, r_1, \dots, r_n)$. (Not all r_k need actually occur in $B^*(p, r_1, \dots, r_n)$, and neither need p .)

The intuitive content of $A_p \leq B_p$ is that propositional letters q different from p occur in B_p in no other context than they occur in A_p . Clearly \leq is transitive. We allow that the sequence $E_1 q, \dots, E_n q$ is empty; this means that $A_p \leq B_p$ implies that if q occurs in B_p , then q occurs in A_p . We have:

2.3 Lemma

- i) Suppose $A_p \leq B_p$ and $A_p \leq C_p$, then $A_p \leq B C_p$.
- ii) Suppose $A_p \leq B(p, p)$, $A_p \leq C_p$ and $A_p \leq D_p$, then $A_p \leq B(C_p, D_p)$.
- iii) Suppose that A_p is of the form $B C_p$, that p really occurs in C_p and that p does not occur in C_q , then $A_p \leq B_p$ and $A_p \leq C_p$.
- iv) If at most the propositional variable p occurs in B_p , then $A_p \leq B A_p$.
- v) Suppose $A(p, q) \leq B(p, q)$, then $A(p, p) \leq B(p, p)$.
- vi) If $A_p = B_p \# C_p$ and p really occurs in A_p , then $A_p \leq B_p$.

Proofs: The proofs of (i) and (ii) are trivial. For (iii), it is sufficient to note that $A^*(p, E_1 q, \dots, E_n q)$ must be of the form $B^*(C^*(p, E_1 q, \dots, E_n q), E_1 q, \dots, E_n q)$. (The occurrence of p in C_p must be real, to

make sure that C_p cannot be a subformula of one of the $E_k q$.) (iv) is easy. Ad (v): suppose $A(p,p)$ is of the form $A^*(p,p,E_1 r, \dots, E_n r)$. This means that $A(p,q)$ is of the form $A^*(p,q,E_1 r, \dots, E_n r)$. So $B(p,q)$ must be of the form $B^*(p,q,E_1 r, \dots, E_n r)$. Clearly q does not occur in the $E_k r$, so the form for $B(p,p)$ we are looking for is $B^*(p,p,E_1 r, \dots, E_n r)$. For (vi), note that $A^*(p,E_1 q, \dots, E_n q)$ must be of the form $B^*(p,E_1 q, \dots, E_n q) \# C^*(p,E_1 q, \dots, E_n q)$. \square

2.4 Explicit Definability Theorem, part 2

Let U be any extension of SR_0 satisfying:

FIX Every formula A_p of the form $\Box B_p$ or $B_p \# C_p$ has a fixed point D such that $A_p \leq D$.

For every formula A_p with p modalized, there is a formula D such that: p does not occur in D , $A_p \leq D$ and $U \vdash D \leftrightarrow AD$.

Proof: Let p be modalized in A_p . Let $A_p = B(C_1 p, \dots, C_n p)$, where the $C_k p$ are either of the form $\Box E_p$ or of the form $E_p \# F_p$ and where p does not occur in $B(q_1, \dots, q_n)$.

Our proof is by induction on n . First suppose $n=1$. Suppose A_p is of the form $B C_p$, where p does not occur in Bq and C_p is either of the form $\Box D_p$ or $D_p \# E_p$. We may assume that p really occurs in C_p . Let D be the fixed point of $C B_p$ guaranteed by **FIX**. We show that $\vdash B D \leftrightarrow A B D$. We have $\vdash D \leftrightarrow C B D$. So by S_1 : $\vdash B D \leftrightarrow B C B D$, and clearly $B C B D = A B D$. Trivially p does not occur in $B D$. We have: $A_p \leq B_p$, $A_p \leq C_p$, hence $A_p \leq C B_p$. Because $C B_p \leq D$, it follows that $A_p \leq D$ and thus $A_p \leq B D$.

For the induction step we have to show how to reduce the number of 'components' in A_p . Suppose q does not occur in A_p . Define $A^*(p,q)$ by $B(C_1 p, \dots, C_{n-1} p, C_n q)$. $A^*(p,q)$ has $n-1$ components in which p occurs, so we may apply the induction hypothesis to get D_q with $A^*(p,q) \leq D_q$ and $\vdash D_q \leftrightarrow A^*(D_q, q)$. Clearly D_q can be written as $F C_n q$, where q does not occur in F . Applying the basis step of our induction to $F C_n p$ we find an E with: $\vdash E \leftrightarrow D E$, and thus $\vdash E \leftrightarrow A^*(D E, E)$. By S_1 it follows that $\vdash E \leftrightarrow A^*(E, E)$. Clearly $A^*(E, E) = A E$. Evidently p does not occur in E . Finally: $A_p = A^*(p,p) \leq D_p \leq E$. \square

2.5 Corollary

(a) For every formula A_p with p modalized, there is a formula D such that p does not occur in D and $SR_1 \vdash D \leftrightarrow A D$.

(b) For every formula A_p in the language of interpretability logic with p modalized, there is a formula D such that p does not occur in D and $ILP \vdash D \leftrightarrow A D$.

Proof: (a) The fixed points D for formulas A_p of the form $\Box B_p$ or $B_p \# C_p$ which SR_1 has by the Explicit Definability Theorem, part 1, are $\Box B \top$ and $B \top \# C \top$ respectively. Since, by lemma 2.3(i) and (iv), $\Box B_p \leq \Box B \top$ and $B_p \# C_p \leq B \top \# C \top$, SR_1 satisfies **FIX**.

(b) Follows immediately from (a). \square

Corollary 2.5(a) is Smoryński's version of the *Explicit Definability Theorem* with a proof along the lines of his "slightly easier proof" (see Smoryński[85], p.81). The original theorem was due to De Jongh and Sambin. Our proof differs only in two minor details from Smoryński's. First, for our purpose of proving the theorem for **IL**, it is essential that 2.4 is not proven in **SR₁**, as **SR₁** is valid for **ILP**, but not for **IL**, or even for **ILM**. Secondly, the artifice of using \leq was added, because the generality of theorem 2.4 forced us to be more explicit than usual about the property of the fixed points needed to get the proof to work. Surely our choice of the property ' $\text{Ap} \leq \text{D}$ ' is not the most parsimonious one, but we submit that it is fairly natural.

3 Explicit fixed points for **IL**

As is easily seen **IL** satisfies the principle E of the system **SR₀**. So, the Uniqueness Theorem, 2.1, holds for **IL**. On the other hand, using **IL**-models, one can show that **IL** does not satisfy L3'. So, the proof of the *Explicit Definability Theorem*, part 1, is not available for **IL**. Thus we have to provide a different proof for *Explicit Definability*, part 1 for **IL**. This is the main aim of this section. Before giving the proof we list some theorems of **IL**.

Define: $A \equiv B := (A \triangleright B) \wedge (B \triangleright A)$.

K1 $\vdash A \equiv (A \vee \diamond A)$ J1, J5, J3

Let $\phi A := (A \vee \diamond A)$, $\psi A := (A \wedge \square \neg A)$, then by L1-L3:

K2 $\vdash \phi A \leftrightarrow \phi \phi A$
 $\vdash \phi A \leftrightarrow \phi \psi A$
 $\vdash \psi A \leftrightarrow \psi \psi A$
 $\vdash \psi A \leftrightarrow \psi \phi A$

Immediate consequences of the above are:

K3 $\vdash A \triangleright A \wedge \square \neg A$

K4 $\vdash A \equiv A \wedge \square \neg A$

Note that: K4 is an alternative for axiom J5.

K5 $\vdash A \triangleright \perp \rightarrow \square \neg A$ J4

Feferman's Principle is the following:

F $\vdash \diamond A \rightarrow \neg(A \triangleright \diamond A)$

F is *not* derivable in **IL**. However, the following weakening of F is derivable:

$$\text{K6} \quad \vdash \Diamond A \triangleright \neg(A \triangleright \Diamond A)$$

Proof: By the above it is sufficient to show: **IL** $\vdash (\Diamond A \wedge \Box \neg \Diamond A) \rightarrow \neg(A \triangleright \Diamond A)$. We have:

$$\begin{aligned} \vdash (\Diamond A \wedge \Box \neg \Diamond A \wedge (A \triangleright \Diamond A)) &\rightarrow (\Diamond A \wedge \Box \neg A \wedge (A \triangleright \Diamond A)) \\ &\rightarrow (\Diamond A \wedge A \triangleright \perp) \\ &\rightarrow (\Diamond A \wedge \Box \neg A) \\ &\rightarrow \perp \quad \square \end{aligned}$$

Start of the proof of Explicit Definability, part 1.

$$\text{E1} \quad \text{Suppose: } \vdash \Box \neg A \triangleright T \rightarrow C, \text{ then } \vdash A \triangleright T \wedge \Box \neg A \triangleright T \leftrightarrow AC \wedge \Box \neg AC.$$

Proof: The " \rightarrow " side is immediate, because $\Box \neg A \triangleright T \rightarrow \Box^+(C \leftrightarrow T)$.

" \leftarrow " Suppose $\vdash \Box \neg A \triangleright T \rightarrow C$. Reason inside the " \vdash ": Suppose AC and $\Box \neg AC$. We have: $\Box(\Box \neg A \triangleright T \rightarrow \Box^+(C \leftrightarrow T))$. Combining this with $\Box \neg AC$ we get: $\Box(\Box \neg A \triangleright T \rightarrow \neg A \triangleright T)$. Hence by Löb's Principle: $\Box \neg A \triangleright T$. It follows that $\Box^+(C \leftrightarrow T)$. Combining this with AC we find A \triangleright T. \square

$$\text{E2} \quad \text{Suppose: } \vdash \Box \neg A \triangleright T \rightarrow C, \text{ then } \vdash A \triangleright T \equiv AC. \quad \text{E1, K4}$$

$$\text{E3} \quad \vdash A \triangleright T \equiv A(A \triangleright T \triangleright B \Box \neg A \triangleright T)$$

Proof: We have $\vdash \Box \neg A \triangleright T \rightarrow A \triangleright T \triangleright B \Box \neg A \triangleright T$. Apply E2. \square

$$\text{E4} \quad \vdash \Box \neg B \Box \neg A \triangleright T \rightarrow (A \triangleright T \triangleright B \Box \neg A \triangleright T \leftrightarrow \Box \neg A \triangleright T)$$

$$\begin{aligned} \text{Proof: } \vdash \Box \neg B \Box \neg A \triangleright T &\rightarrow (A \triangleright T \triangleright B \Box \neg A \triangleright T \leftrightarrow A \Box \neg A \triangleright T \triangleright \perp) \\ &\leftrightarrow \Box \neg A \triangleright T \quad \square \end{aligned}$$

$$\text{E5} \quad \vdash \Box \neg B \Box \neg A \triangleright T \rightarrow \Box^+(A \triangleright T \triangleright B \Box \neg A \triangleright T \leftrightarrow \Box \neg A \triangleright T)$$

$$\text{E6} \quad \vdash B \Box \neg A \triangleright T \wedge \Box \neg B \Box \neg A \triangleright T \leftrightarrow B(A \triangleright T \triangleright B \Box \neg A \triangleright T) \wedge \Box \neg B(A \triangleright T \triangleright B \Box \neg A \triangleright T)$$

Proof: " \rightarrow ": immediate by E5 and S_2 . For the " \leftarrow "-side it is clearly sufficient to show:

$$\vdash \Box \neg B(A \triangleright T \triangleright B \Box \neg A \triangleright T) \rightarrow \Box \neg B \Box \neg A \triangleright T$$

This follows by:

$$\begin{aligned} \vdash \Box \neg B(A \triangleright T \triangleright B \Box \neg A \triangleright T) &\rightarrow \Box(\Box \neg B \Box \neg A \triangleright T \rightarrow \neg B \Box \neg A \triangleright T) \quad (\text{E5, } S_2) \\ &\rightarrow \Box \neg B \Box \neg A \triangleright T \quad \square \end{aligned}$$

$$\text{E7} \quad \vdash B \Box \neg A \triangleright T \equiv B(A \triangleright T \triangleright B \Box \neg A \triangleright T) \quad \text{E6, K4}$$

$$\text{E8} \quad \vdash A \triangleright T \triangleright B \Box \neg A \triangleright T \leftrightarrow A(A \triangleright T \triangleright B \Box \neg A \triangleright T) \triangleright B(A \triangleright T \triangleright B \Box \neg A \triangleright T) \quad \text{E3, E7}$$

End of the proof of Explicit Definability, part 1.

It is easy to see that p does not occur in $A \top \triangleright B \Box \neg A \top$. We have: $(A \triangleright B) \leq (A \top \triangleright B \Box \neg A \top)$. For assume that p really occurs in $A \triangleright B$. By 2.3: $(A \triangleright B) \leq A \top \leq \Box \neg A \top$. Also $(A \triangleright B) \leq \top$. Combining by 2.3(ii) we find: $(A \triangleright B) \leq (A \top \triangleright B \Box \neg A \top)$. So, we can apply 2.4 and conclude Explicit Definability for **IL**:

for every formula A with p modalized, there is a formula D such that:

p does not occur in D , and $\mathbf{IL} \vdash D \leftrightarrow A$.

4 The system **ILW**

The principle **W** is very powerful. It can be viewed (in our limited context) as a generalization both of Gödel's Second Incompleteness Theorem and of Gödel's Completeness Theorem (in the guise of the Interpretation Existence Lemma). To illustrate this we show that **ILW** can be axiomatized as follows:

- L1 $\vdash A \Rightarrow \vdash \Box A$
- L2 $\vdash \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$
- J1 $\vdash \Box(A \rightarrow B) \rightarrow A \triangleright B$
- J2 $\vdash (A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C$
- J3 $\vdash (A \triangleright C) \wedge (B \triangleright C) \rightarrow A \vee B \triangleright C$
- J4 $\vdash A \triangleright B \rightarrow (\Diamond A \rightarrow \Diamond B)$
- W $\vdash A \triangleright B \rightarrow A \triangleright B \wedge \Box \neg A$

First prove Feferman's principle **F** by substituting $\Diamond A$ for B in **W** (this uses L1, L2, J1, J2). Löb's Principle (**L4**) then follows from **F**:

$$\begin{aligned} \vdash \Box(\Box A \rightarrow A) &\rightarrow \Box(\neg A \rightarrow \Diamond \neg A) \\ &\rightarrow \neg A \triangleright \Diamond \neg A \\ &\rightarrow \neg \Diamond \neg A \\ &\rightarrow \Box A \end{aligned}$$

Using **L4** one derives **L3** by a well-known trick. Next we derive **K2**. Using **K2** and $\vdash A \equiv A \wedge \Box \neg A$ which is immediate by **W**, we get: $\vdash A \equiv A \vee \Diamond A$ and hence, by J1, J5.

W is not derivable in **IL**. To show this we need some model theory: we use Frank Veltman's **IL**-models. An **IL**-model M is of the form: $\langle K, R, S, \Vdash \rangle$, where: K is non-empty; R is a binary relation on K , which is transitive, upwards well-founded; S is a ternary relation on K , which we treat as a K -indexed set of binary relations S_k on K ; the S_k are reflexive, transitive; we have: $kRmS_k n \Rightarrow kRn$ and $kRmRn \Rightarrow mS_k n$; \Vdash is a forcing relation on M , where R is the accessibility relation for \Box and:

$$k \Vdash A \triangleright B \Leftrightarrow \text{for all } m \text{ with } kRm \text{ and } m \Vdash A \text{ there is an } n \text{ with } mS_k n \text{ and } n \Vdash B.$$

It is easy to show that **IL** is valid in **IL**-models, and **IL** is complete w.r.t. (finite) **IL**-models (De Jongh & Veltman[88]).

Consider the **IL**-model on $\{\alpha, \beta, \gamma\}$ generated by $\alpha R \beta R \gamma$, $\gamma S_{\alpha} \beta$, $\gamma \Vdash p$. Clearly $\alpha \Vdash p \triangleright \diamond p$, but $\alpha \not\Vdash \Box \neg p$. Hence Feferman's Principle doesn't hold at α and so a fortiori **W** fails.

We show that the Fixed Point of $A p \triangleright B p$ found in Section 3 simplifies in **ILW** to $A \top \triangleright B \top$:

$$\vdash A \top \triangleright B \top \leftrightarrow A \top \triangleright B \Box \neg A \top.$$

Proof:

$$\begin{aligned} \vdash A \top \triangleright B \top &\leftrightarrow A \top \triangleright B \top \wedge \Box \neg A \top \\ &\leftrightarrow A \top \triangleright B \Box \neg A \top \wedge \Box \neg A \top \\ &\leftrightarrow A \top \triangleright B \Box \neg A \top \quad \square \end{aligned}$$

Finally we show that the simplified fixed point doesn't work in **IL**. Consider $q \triangleright \neg p$. The **ILW**-style fixed point in p for this formula is: $q \triangleright \neg \top$, i.e. modulo **IL** provable equivalence: $\Box \neg q$. If this were a fixed point in **IL**, we would have: $\mathbf{IL} \vdash \Box \neg q \leftrightarrow q \triangleright \diamond q$. We have already seen that this is not the case.

References:

- Boolos, G., 1979, *The Unprovability of Consistency*, CUP, London.
- Buss, S., 1985, *Bounded Arithmetic*, Thesis, Princeton University, Princeton. Reprinted: 1986, Bibliopolis, Napoli.
- De Jongh, D.H.J. & Veltman, F., 1988, *Provability Logics for Relative Interpretability*. To appear in the Proceedings of the Heyting Conference, Chaika, Bulgaria, 1988.
- Hájek, P., 1981, *Interpretability in Theories containing Arithmetic II*, *Commentationes Mathematicae Universitatis Carolinae* 22, 667-688.
- Smoryński, C., 1985, *Self-Reference and Modal Logic*, Springer Verlag.
- Švejdar, V., 1983, *Modal Analysis of Generalized Rosser Sentences*, *JSL* 48, 986-999.
- Visser, A., 1988, *Interpretability Logic*, Logic Group Preprint Series nr 40, Dept. of Philosophy, University of Utrecht, Heidelberglaan 2, 3584CS Utrecht. To appear in the Proceedings of the Heyting Conference, Chaika, Bulgaria, 1988.
- Visser, A., 1988P, *Preliminary Notes on Interpretability Logic*, Logic Group Preprint Series nr 29, Dept. of Philosophy, University of Utrecht, Heidelberglaan 2, 3584CS Utrecht.

Logic Group Preprint Series

Department of Philosophy
University of Utrecht
Heidelberglaan 2
3584 CS Utrecht
The Netherlands

- nr. 1 C.P.J. Koymans, J.L.M. Vrancken, *Extending Process Algebra with the empty process*, September 1985.
- nr. 2 J.A. Bergstra, *A process creation mechanism in Process Algebra*, September 1985.
- nr. 3 J.A. Bergstra, *Put and get, primitives for synchronous unreliable message passing*, October 1985.
- nr. 4 A. Visser, *Evaluation, provably deductive equivalence in Heyting's arithmetic of substitution instances of propositional formulas*, November 1985.
- nr. 5 G.R. Renardel de Lavalette, *Interpolation in a fragment of intuitionistic propositional logic*, January 1986.
- nr. 6 C.P.J. Koymans, J.C. Mulder, *A modular approach to protocol verification using Process Algebra*, April 1986.
- nr. 7 D. van Dalen, F.J. de Vries, *Intuitionistic free abelian groups*, April 1986.
- nr. 8 F. Voorbraak, *A simplification of the completeness proofs for Guaspari and Solovay's R*, May 1986.
- nr. 9 H.B.M. Jonkers, C.P.J. Koymans & G.R. Renardel de Lavalette, *A semantic framework for the COLD-family of languages*, May 1986.
- nr. 10 G.R. Renardel de Lavalette, *Strictheidsanalyse*, May 1986.
- nr. 11 A. Visser, *Kunnen wij elke machine verslaan? Beschouwingen rondom Lucas' argument*, July 1986.
- nr. 12 E.C.W. Krabbe, *Naess's dichotomy of tenability and relevance*, June 1986.
- nr. 13 Hans van Ditmarsch, *Abstractie in wiskunde, expertsystemen en argumentatie*, Augustus 1986
- nr. 14 A. Visser, *Peano's Smart Children, a provability logical study of systems with built-in consistency*, October 1986.
- nr. 15 G.R. Renardel de Lavalette, *Interpolation in natural fragments of intuitionistic propositional logic*, October 1986.
- nr. 16 J.A. Bergstra, *Module Algebra for relational specifications*, November 1986.
- nr. 17 F.P.J.M. Voorbraak, *Tensed Intuitionistic Logic*, January 1987.
- nr. 18 J.A. Bergstra, J. Tiuryn, *Process Algebra semantics for queues*, January 1987.
- nr. 19 F.J. de Vries, *A functional program for the fast Fourier transform*, March 1987.
- nr. 20 A. Visser, *A course in bimodal provability logic*, May 1987.
- nr. 21 F.P.J.M. Voorbraak, *The logic of actual obligation, an alternative approach to deontic logic*, May 1987.
- nr. 22 E.C.W. Krabbe, *Creative reasoning in formal discussion*, June 1987.
- nr. 23 F.J. de Vries, *A functional program for Gaussian elimination*, September 1987.
- nr. 24 G.R. Renardel de Lavalette, *Interpolation in fragments of intuitionistic propositional logic*, October 1987. (revised version of no. 15)
- nr. 25 F.J. de Vries, *Applications of constructive logic to sheaf constructions in toposes*, October 1987.
- nr. 26 F.P.J.M. Voorbraak, *Redeneren met onzekerheid in expertsystemen*, November 1987.
- nr. 27 P.H. Rodenburg, D.J. Hoekzema, *Specification of the fast Fourier transform algorithm as a term rewriting system*, December 1987.

- nr. 28 D. van Dalen, *The war of the frogs and the mice, or the crisis of the Mathematische Annalen*, December 1987.
- nr. 29 A. Visser, *Preliminary Notes on Interpretability Logic*, January 1988.
- nr. 30 D.J. Hoekzema, P.H. Rodenburg, *Gauß elimination as a term rewriting system*, January 1988.
- nr. 31 C. Smorynski, *Hilbert's Programme*, January 1988.
- nr. 32 G.R. Renardel de Lavalette, *Modularisation, Parameterisation, Interpolation*, January 1988.
- nr. 33 G.R. Renardel de Lavalette, *Strictness analysis for POLYREC, a language with polymorphic and recursive types*, March 1988.
- nr. 34 A. Visser, *A Descending Hierarchy of Reflection Principles*, April 1988.
- nr. 35 F.P.J.M. Voorbraak, *A computationally efficient approximation of Dempster-Shafer theory*, April 1988.
- nr. 36 C. Smorynski, *Arithmetic Analogues of McAloon's Unique Rosser Sentences*, April 1988.
- nr. 37 P.H. Rodenburg, F.J. van der Linden, *Manufacturing a cartesian closed category with exactly two objects*, May 1988.
- nr. 38 P.H. Rodenburg, J. L.M. Vrancken, *Parallel object-oriented term rewriting : The Booleans*, July 1988.
- nr. 39 D. de Jongh, L. Hendriks, G.R. Renardel de Lavalette, *Computations in fragments of intuitionistic propositional logic*, July 1988.
- nr. 40 A. Visser, *Interpretability Logic*, September 1988.
- nr. 41 M. Doorman, *The existence property in the presence of function symbols*, October 1988.
- nr. 42 F. Voorbraak, *On the justification of Dempster's rule of combination*, December 1988.
- nr. 43 A. Visser, *An inside view of EXP, or: The closed fragment of the provability logic of $I\Delta_0 + \Omega_1$* , February 1989.
- nr. 44 D.H.J. de Jongh & A. Visser, *Explicit Fixed Points in Interpretability Logic*, March 1989.