

Propositional combinations
of Σ -sentences
in Heyting's Arithmetic

A. Visser

θ

π



©1994, Department of Philosophy - Utrecht University

ISBN 90-393-0713-x

ISSN 0929-0710

Dr. A. Visser, Editor

Propositional combinations of Σ -sentences in Heyting's Arithmetic

*Explorations between
intuitionistic propositional logic
and intuitionistic arithmetic*

Albert Visser
Version 1.0

ABSTRACT: This paper is mainly about Boolean combinations of Σ -formulas ($B(\Sigma)$ -formulas) in HA. We prove a theorem guaranteeing that if a $B(\Sigma)$ -formula is provable in HA, then a simpler one is also provable. Our theorem yields a characterization of the derived rules of HA for Σ -substitutions. Another application is a decision procedure for the closed fragment of the provability logic of HA. The proof of our theorem leads us to such varied subjects as NNIL-formulas, robust formulas and preservativity notions.

1 Introduction

1.1 The problem and its solution: This paper is about the simple question:

What more do we know, when we know that a Boolean (or, if you prefer, a Brouwerian) combination of Σ -sentences is provable in Heyting's Arithmetic?

The answer takes the form:

We often know that a better Boolean combination of the same Σ -sentences is provable in Heyting's Arithmetic. Moreover the better Boolean combination can be found independently of the specific Σ -sentences under consideration.

Here *better* has two components. It means 'stronger' in the sense that the stronger sentence implies the weaker one in the Intuitionistic Propositional Calculus (IPC). It also means 'simpler' in the sense that the simpler formula lies in a specific class NNIL, the class of propositional formulas with No Nestings of Implications to the

Left. The NNIL-formulas are modulo IPC-provable equivalence precisely the formulas preserved under sub-models in the usual Kripke semantics for IPC. The qualification *often* in the answer is just a gloss for:

unless our Boolean combination is already in NNIL.

We give a more formal statement of our result. Let A be a propositional formula. We let \bar{f} range over assignments of Σ -sentences to (at least) the propositional formulas occurring in A . We write ' $A[\bar{f}]$ ' for the result of substituting $\bar{f}(p)$ for p in A . Our theorem says: there is an $A^* \in \text{NNIL}$ such that:

- a) $\text{IPC} \vdash A^* \rightarrow A$, and
- b) for all \bar{f} : $\text{HA} \vdash A[\bar{f}] \Rightarrow \text{HA} \vdash A^*[\bar{f}]$.

A better answer to our question is not possible, at least if we abstract away from the specific Σ -sentences substituted. We will show:

- c) For every propositional formula B such that $\text{IPC} \not\vdash A^* \rightarrow B$, there is an \bar{f} such that:
 $\text{HA} \vdash A^*[\bar{f}]$ and $\text{HA} \not\vdash B[\bar{f}]$.

Or contraposed:

- d) For every propositional formula B :
 $\forall \bar{f} (\text{HA} \vdash A^*[\bar{f}] \Rightarrow \text{HA} \vdash B[\bar{f}]) \Rightarrow \text{IPC} \vdash A^* \rightarrow B$.

From (a),(b) we get:

- e) $\forall \bar{f} (\text{HA} \vdash A[\bar{f}] \Leftrightarrow \text{HA} \vdash A^*[\bar{f}])$.

And hence (d) is equivalent with:

- f) For every propositional formula B :
 $\forall \bar{f} (\text{HA} \vdash A[\bar{f}] \Rightarrow \text{HA} \vdash B[\bar{f}]) \Rightarrow \text{IPC} \vdash A^* \rightarrow B$.

So A^* is the strongest 'general' improvement of A . Of course, for specific \bar{f} often better improvements are possible, but A^* is the best one that works for all \bar{f} . Note that (f) implies (a) by substituting A for B . So our results are summarized by (b) and (f).

1.2 What more is there in the paper? Our main result has two immediate applications. Firstly it yields a description of the derived rules of HA for Σ -substitutions. Secondly we can read off a decision procedure for the closed fragment of the provability logic of HA. Thus we solve a variant of Friedman's 35th problem (see Friedman[75]). (This solution occurs already in Visser[85])

Both for the smooth formulation of our main result and for its proof we have to develop some further material. Part of it involves the theory of certain classes of propositional formulas, like NNIL and ROB (to be introduced in the paper). A second part is the study of certain relations between formulas, called semi-consequence relations. The semi-consequence relations include derived rules and preservativity rela-

tions (dual to conservativity relations). One preservativity relation, namely Σ -preservativity (the ‘dual’ of Π -conservativity), will be studied in some detail in a long excursion (section 7).

The main technical result of the paper is the specification and verification of the NNIL-algorithm in section 5. Our results in both propositional logic and arithmetic are immediate consequences of the existence of this algorithm.

The contents of the paper is as follows:

Section 2	preliminaries to propositional logic
Section 3	consequence relations and preservativity notions
Section 4	robust formulas
Section 5	the NNIL-algorithm
Section 6	basic facts and notations for arithmetic
Section 7	digression: assorted facts on Σ -preservativity
Section 8	closure properties of Σ -preservativity over HA
Section 9	on Σ -substitutions
Section 10	the closed fragment of the provability logic of HA

1.3 Environment and history of the paper: This paper is part of the wider study of substitutions considered as a kind of semantics. We point to some examples of related work. First there is, of course, the project of provability and interpretability logic. See e.g. Smoryński[85], Boolos[93], Berarducci[90], Visser[90], EA[91], Shavrukov[93], Zambella[94]. Then there is the work on derived rules by Rybakov. See e.g. Rybakov[92]. For a systematical exploration of the notion of derived rule, see FHY[92]. Finally there is the most immediate environment: research on substitutions in intuitionistic logic of propositional formulas or arithmetical sentences. See e.g. De Jongh[82], Smoryński[73], Leivant[75,80,81], Van Oosten[91], DV[94], DC[?].

The present paper is, partly, a remake of the unpublished paper *Evaluation, provably deductive equivalence in Heyting's Arithmetic* (Visser[85]). The work in Visser[85] was inspired by Dick de Jongh's results on propositional formulas of one variable, reported in De Jongh[82]. Since Visser[85] was written, the development of interpretability logic took place (see e.g. Berarducci[90] and Visser[90]). Interpretability logic inspired me to give the notion of preservativity a more prominent role in the presentation.

In Visser[85] we proved the identity of certain formula classes NNIL and ROB. Johan van Benthem gave, independently, around 1984 an alternative proof of the fact that

NNIL=ROB (see his Van Benthem[91]). Gerard Renardel proved (also around 1984, see his Renardel[86]), that NNIL satisfies both left and right interpolation. The work of both van Benthem and Renardel is contained and extended in DRVV[94]. This last paper should be considered as a companion paper of the present paper.

Dick de Jongh and Albert Visser recently reopened the research on Heyting algebras of arithmetical theories. A number of concerns of the present paper (derived rules, exact formulas) reappear in their paper DV[94], which is another companion paper of the present paper.

1.4 Acknowledgements: In various stages of research I benefited from the work, the wisdom and/or the advice of: Johan van Benthem, Dirk van Dalen, Dick de Jongh, Karst Koymans, Piet Rodenburg, Volodya Shavrukov, Rick Statman, Anne Troelstra and Domenico Zambella.

1.5 Prerequisites: Some knowledge of Troelstra[73] or TV[88a,b] is certainly beneficial. At some places I will make use of results from Visser[82] and DV[94].

2 Preliminaries to propositional logic: $\mathfrak{P}, \mathfrak{Q}, \dots$ will be *sets of propositional variables*. $\mathbf{p}, \mathbf{q}, \mathbf{r}, \dots$ will be *finite sets of propositional variables*. We define $\mathfrak{L}(\mathfrak{P})$ as the smallest set \mathfrak{S} such that:

- $\mathfrak{P} \subseteq \mathfrak{S}$, $\top, \perp \in \mathfrak{S}$,
- If $A, B \in \mathfrak{S}$, then $(A \wedge B), (A \vee B), (A \rightarrow B) \in \mathfrak{S}$.

SUB(A) is the set of subformulas of A. By convention we will count \perp a subformula of any A. PV(A) is the set of propositional variables occurring in A.

We suppose that the reader is familiar with Kripke models for IPC (see TV[88a], or Smoryński [73]). Our treatment here is mainly to fix notations. A model is a structure $\mathbb{K} = \langle K, \leq, \models, \mathfrak{P} \rangle$, where K is a non-empty set of nodes, \leq is a partial ordering. \models is the atomic forcing relation for \mathfrak{P} : it is a relation between nodes and propositional atoms in \mathfrak{P} , satisfying: $k \leq k'$ and $k \models p \Rightarrow k' \models p$ (persistence). \mathbb{K} is a *\mathfrak{P} -model* if $\mathfrak{P}_{\mathbb{K}} = \mathfrak{P}$. $\text{Mod}(\mathfrak{P})$ is the set of \mathfrak{P} -models.

Consider $\mathbb{K} = \langle K, \leq, \models, \mathfrak{P} \rangle$. We define $k \models A$ for $A \in \mathfrak{L}(\mathfrak{P})$ in the standard way. We write $\mathbb{K} \models A$ for: $\forall k \in K \ k \models A$.

A model \mathbb{K} is *finite* if both $K_{\mathbb{K}}$ and $\mathfrak{P}_{\mathbb{K}}$ are finite.

A *rooted model* \mathbb{K} is a structure $\langle K, \leq, \models, \mathfrak{P} \rangle$, where $\langle K, \leq, \models, \mathfrak{P} \rangle$ is a model and where $\hat{b} \in K$ is the bottom element w.r.t. \leq . The set of rooted models is $R\text{mod}$.

For any $k \in K$ $\mathbb{K}[k]$ is the model $\langle K', k, \leq', \models', \mathfrak{P} \rangle$, where $K' := \uparrow k := \{k' \mid k \leq k'\}$ and where \leq' and \models' are the restrictions of \leq respectively \models to K' . (We will often simply write \leq and \models for \leq' and \models' .)

2.1 The Henkin construction: A set $X \subseteq \mathcal{L}(\mathfrak{P})$ is *adequate* if it is closed under subformulas and contains \perp . A set Γ is X -saturated if:

- (i) $\Gamma \subseteq X$, (ii) $\Gamma \not\models \perp$, (iii) $\Gamma \vdash A, A \in X \Rightarrow A \in \Gamma$,
- (iv) $\Gamma \vdash (B \vee C), (B \vee C) \in X \Rightarrow B \in \Gamma$ or $C \in \Gamma$.

The Henkin model $\mathbb{H}_X(\mathfrak{P})$ is the \mathfrak{P} -model with as nodes the X -saturated sets Δ and as accessibility relation \subseteq . The atomic forcing in the nodes is given by: $\Gamma \models p \Leftrightarrow p \in \Gamma$. We have by a standard argument: for $A \in X$: $\Gamma \models A \Leftrightarrow A \in \Gamma$. Note that if X is finite, then $\mathbb{H}_X(\mathfrak{P})$ is finite. A direct consequence of the Henkin construction is the Kripke Completeness Theorem. Let $\mathfrak{P} \supseteq \text{PV}(A)$, then:

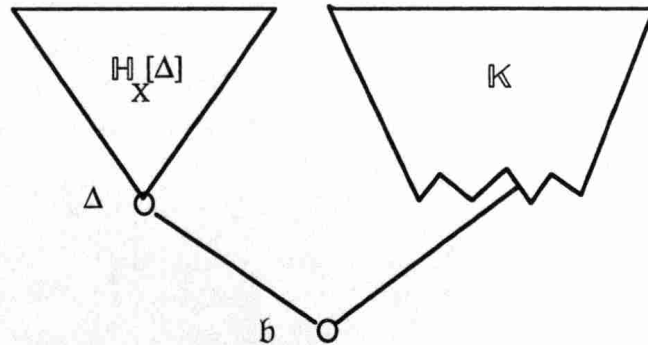
$$\text{IPC} \vdash A \Leftrightarrow \text{for all (finite) } \mathfrak{P}\text{-models } \mathbb{K}: \mathbb{K} \models A.$$

2.2 Further Definitions

- i) Let \mathbf{K} be a set of \mathfrak{P} -models. $M(\mathbf{K})$ is the \mathfrak{P} -model with nodes $\langle k, \mathbb{K} \rangle$ for $k \in K_{\mathbb{K}}$, $\mathbb{K} \in \mathbf{K}$ and ordering: $\langle k, \mathbb{K} \rangle \leq \langle m, \mathbb{M} \rangle \Leftrightarrow \mathbb{K} = \mathbb{M}$ and $k \leq_{\mathbb{K}} m$. As atomic forcing we define: $\langle k, \mathbb{K} \rangle \models p \Leftrightarrow k \models_{\mathbb{K}} p$. (In practice we will forget the second components of the new nodes, pretending the domains to be already disjoint.)
- ii) Let \mathbb{K} be a \mathfrak{P} -model. $B(\mathbb{K})$ is the rooted \mathfrak{P} -model obtained by adding a new bottom \hat{b} to \mathbb{K} and by taking: $\hat{b} \models p \Leftrightarrow \mathbb{K} \models p$. We put $\text{Glue}(\mathbf{K}) := B(M(\mathbf{K}))$. \circ

We will assume below that \mathfrak{P} is fixed. We will often notationally suppress it.

2.3 Push Down Lemma: Let X be adequate. Suppose Δ is X -saturated and $\mathbb{K} \models \Delta$. Then $\text{Glue}(\mathbb{H}_X[\Delta], \mathbb{K}) \models \Delta$.



Proof: We show by induction on $A \in X$ that $b \models A \Leftrightarrow A \in \Delta$. The cases of atoms, conjunction and disjunction are trivial. If $(B \rightarrow C) \in X$ and $b \models (B \rightarrow C)$, then $\Delta \models (B \rightarrow C)$ and hence $(B \rightarrow C) \in \Delta$. Conversely suppose $(B \rightarrow C) \in \Delta$. If $b \not\models B$, we are easily done. If $b \models B$, then, by the Induction Hypothesis: $B \in \Delta$, hence $C \in \Delta$ and, by the Induction Hypothesis: $b \models C$. \square

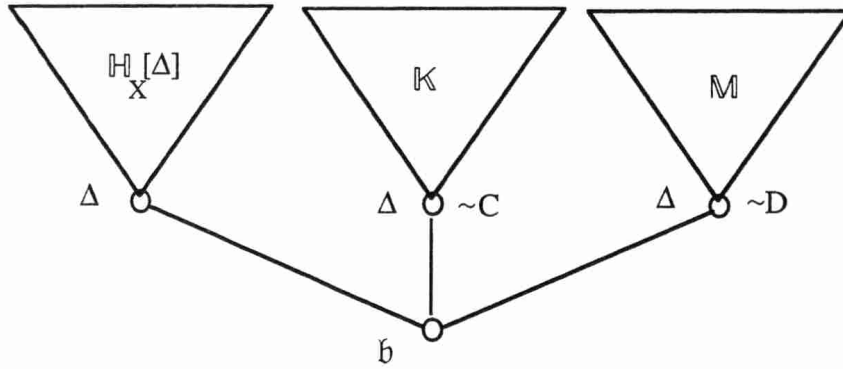
We say that Δ is (\mathfrak{P}) -prime if it is consistent and:

for every $(C \vee D) \in \mathfrak{L}(\mathfrak{P})$: $\Delta \vdash (C \vee D) \Rightarrow \Delta \vdash C$ or $\Delta \vdash D$.

A formula A is prime if $\{A\}$ is prime.

2.4 Theorem: Suppose X is adequate and Δ is X -saturated. Then Δ is prime.

Proof: Δ is consistent by definition. Suppose $\Delta \vdash C \vee D$ and $\Delta \not\vdash C$ and $\Delta \not\vdash D$. Suppose $K \models \Delta$, $K \not\models C$, $M \models \Delta$ and $M \not\models D$. Consider $\text{Glue}(\mathbb{H}_X(\Delta), K, M)$. By 2.3 we have: $b \models \Delta$. On the other hand, by persistence: $b \not\models C$ and $b \not\models D$. Contradiction. \square



($\sim C$ exhibited next to a node means that C is not forced; this is not to be confused with $\neg C$ exhibited next to a node, which means that $\neg C$ is forced.)

2.5 Theorem: Consider any formula A . The formula A can be written (modulo IPC-provable equivalence) as a disjunction of prime formulas C . Moreover these C are conjunctions of implications and propositional variables in $\text{SUB}(A)$.

Proof: Consider a $\text{SUB}(A)$ -saturated Δ . Let $\text{IP}(\Delta)$ be the set of implications and atoms of Δ . It is easily seen that $\text{IPC} \vdash \bigwedge \text{IP}(\Delta) \Leftrightarrow \bigwedge \Delta$. Take:

$$D := \bigvee \{ \bigwedge \text{IP}(\Delta) \mid \Delta \text{ is } \text{SUB}(A)\text{-saturated and } A \in \Delta \}.$$

Trivially: $\text{IPL} \vdash D \rightarrow A$. On the other hand if $\text{IPC} \not\vdash A \rightarrow D$, then by a standard construction there is a $\text{SUB}(A)$ -saturated set Γ such that $A \in \Gamma$ and $\Gamma \not\models D$. Quod non. \square

2.5.1 Remark: Note that in the definition of D in 2.5, we can restrict ourselves to the \subseteq -minimal $\text{sub}(A)$ -saturated Δ with $A \in \Delta$.

We define a measure of complexity ρ , which counts the left-nesting of \rightarrow , as follows:

- $\rho(p) := \rho(\perp) := \rho(\top) := 0$
- $\rho(A \wedge B) := \rho(A \vee B) := \max(\rho(A), \rho(B))$
- $\rho(A \rightarrow B) := \max(\rho(A) + 1, \rho(B))$.

$\text{NNIL}(\mathfrak{B}) := \{A \in \mathfrak{L}(\mathfrak{B}) \mid \rho(A) \leq 1\}$. In other words NNIL is the class of formulas without nestings of implications to the left. An example of as NNIL-formula is:

$$(p \rightarrow (q \vee (s \rightarrow t))) \wedge ((q \vee r) \rightarrow s).$$

It is easy to see that modulo IPC-provable equivalence each NNIL-formula can be rewritten to a NNIL_0 -formula, i.e. a formula in which in front of implications only single atoms occur. For more information about NNIL, see DV[94] and DRVV[94].

3 Consequence relations & preservativity notions: An important tool of the present paper is the use of *semi-consequence relations* (defined below) and *preservativity relations* (defined below). ' \triangleright ' will range over semi-consequence relations. Let \vdash stand for derivability in IPC.

Let \mathfrak{B} be a language (for propositional or predicate logic) and let T be a theory in \mathfrak{B} . A *semi-consequence relation on \mathfrak{B} over T* is a binary relation on \mathfrak{B} satisfying:

- A1 $A \vdash_T B \Rightarrow A \triangleright B$
- A2 $A \triangleright B \text{ and } B \triangleright C \Rightarrow A \triangleright C$
- A3 $C \triangleright A \text{ and } C \triangleright B \Rightarrow C \triangleright (A \wedge B)$.

The name 'semi-consequence relation' is ad hoc in this paper. We take:

- $A \equiv B :\Leftrightarrow A \triangleright B \text{ and } B \triangleright A$.

If we don't specify the theory with the semi-consequence relation, it's always supposed to be over IPC.

A further salient principle is:

- B1 $A \triangleright C \text{ and } B \triangleright C \Rightarrow (A \vee B) \triangleright C$.

A relation satisfying A1-A3, B1 is called a *nearly-consequence relation*. Note that \vdash_T is a nearly-consequence relation over T .

3.1 Conventions: We will use $X \vdash Y$, $X \triangleright Y$ for respectively $\bigwedge X \vdash \bigvee Y$ and $\bigwedge X \triangleright \bigvee Y$, where X and Y are finite sets of formulas. ($\bigwedge \emptyset := \top$, $\bigvee \emptyset := \perp$). We treat implications similarly.

We write X, Y for: $X \cup Y$, X, A for $X \cup \{A\}$, etc. \circ

Nearly-consequence relations over T can be alternatively described in Genzen style as follows. Nearly-consequence relations are consequence relations satisfying:

$$A1' \quad X \vdash_T Y \Rightarrow X \triangleright Y$$

$$\text{Thin} \quad X \triangleright Y \Rightarrow X, Z \triangleright Y, U$$

$$\text{Cut} \quad X \triangleright Y, A \text{ and } Z, A \triangleright U \Rightarrow X, Z \triangleright Y, U$$

We take the permutation rules to be implicit in the set notation. We leave it to the reader to check the equivalence of the Genzen style principles with A1-A3, B1.

3.2 Fact: Let \triangleright be a semi-consequence relation on $\mathcal{L}(\mathfrak{B})$. Suppose there is a $\Phi: \mathcal{L}(\mathfrak{B}) \rightarrow \mathcal{L}(\mathfrak{B})$, such that for all $A, B \in \mathcal{L}(\mathfrak{B})$: $A \triangleright B \Leftrightarrow \Phi(A) \vdash B$. Then:

- i) If for all $B \in \mathcal{L}(\mathfrak{B})$: $A \triangleright B \Leftrightarrow C \vdash B$, then $\vdash C \leftrightarrow \Phi(A)$.
- ii) Φ considered as a function on the Heyting algebra of IPC is a co-closure operation, i.e.:
 - $\Phi(A) \vdash A$ (co-inductivity)
 - $A \vdash B \Rightarrow \Phi(A) \vdash \Phi(B)$ (monotonicity)
 - $\vdash \Phi(\Phi(A)) \leftrightarrow \Phi(A)$ (idempotency).

In fact we have more than monotonicity, viz:

$$\bullet \quad A \triangleright B \Leftrightarrow \Phi(A) \vdash \Phi(B).$$

- iii) Suppose \triangleright is a nearly-consequence relation. Then $\vdash \Phi(A \vee B) \leftrightarrow (\Phi(A) \vee \Phi(B))$.

Proof: The simple proofs are left to the reader. \square

Some principles involving implication play an important role in the paper. To introduce them we need a syntactical operation. We define the operation $.$ on propositional formulas as follows:

- $[B]p := p$, $[B]\top := \top$, $[B]\perp := \perp$,
- $.$ commutes with \wedge and \vee ,
- $[B](C \rightarrow D) := (B \rightarrow (C \rightarrow D))$.

Note that $.$ does not preserve provable equivalence in the second component. Note also that: $\vdash B \rightarrow ([B]C \leftrightarrow C)$. Define: $[B]X := \{[B]C \mid C \in X\}$.

We have the following principles for implication for semi-consequence relations on $\mathcal{L}(\mathfrak{B})$:

$$B2 \quad A \triangleright B \Rightarrow (p \rightarrow A) \triangleright (p \rightarrow B)$$

$$B3 \quad \text{Let } X \text{ be a finite set of implications and let } Y := \{C \mid (C \rightarrow D) \in X\} \cup \{B\}. \text{ Take } A := \bigwedge X \text{ then: } (A \rightarrow B) \triangleright [A]Y$$

We give an example of an instance of B3. Let e.g.

$$A := ((p \rightarrow q) \wedge (r \rightarrow (s \vee t))), C := (A \rightarrow (u \vee (p \rightarrow r))),$$

then we have:

$$C \triangleright ([A]p \vee [A]r \vee [A](u \vee (p \rightarrow r))), \text{ i.e.:}$$

$$C \triangleright (p \vee r \vee u \vee (A \rightarrow (p \rightarrow r))).$$

3.3 Open problem: Is it possible to axiomatize B3 by a finite number of traditional schemes? We do not allow finite sets or syntactical operations to occur in such schemes. However we *do* allow schematic letters ranging over arbitrary formulas and schematic letters ranging over propositional variables. (So e.g. B2 is acceptable as a scheme.) \circ

We say that a relation satisfying A1-A3, B1-B3 is a σ -relation.

Consider again any language \mathfrak{B} of propositional or predicate logic. Let $X \subseteq \mathfrak{B}$ and let T be any theory in that language. Let \mathfrak{F} be a set of functions from \mathfrak{P} to \mathfrak{B} . We write:

- For $A, B \in \mathfrak{B}$: $A \triangleright_{T, X} B \Leftrightarrow \forall C \in X (C \vdash_T A \Rightarrow C \vdash_T B)$,
- Let $A \in \mathfrak{L}(\mathfrak{P})$. $A[\mathfrak{f}]$ is the result of substituting $\mathfrak{f}(p)$ for p in A for each $p \in \mathfrak{P}$.
- For $A, B \in \mathfrak{L}(\mathfrak{P})$: $A \triangleright_{T, X, \mathfrak{F}} B \Leftrightarrow \forall \mathfrak{f} \in \mathfrak{F} A[\mathfrak{f}] \triangleright_{T, X} B[\mathfrak{f}]$.

If $Y \subseteq \mathfrak{B}$ and $\mathfrak{F} = Y^{\mathfrak{P}}$, we will write $\triangleright_{T, X, Y} B$ for $\triangleright_{T, X, \mathfrak{F}} B$. If X or \mathfrak{F} are singletons we will omit the singleton brackets. If $\mathfrak{B} = \mathfrak{L}(\mathfrak{P})$ and $T = \text{IPC}$, we will often omit 'T' in the index. We will call $\triangleright_{T, X, \mathfrak{F}}$ T, X, \mathfrak{F} -preservativity, etcetera. We will call the $\triangleright_{T, X, \mathfrak{F}}$ and the $\triangleright_{T, X}$ *preservativity relations*. We will call the $\triangleright_{T, X}$ *pure preservativity relations*.

Clearly $\triangleright_{T, X}$ is a semi-consequence relation over T and $\triangleright_{T, X, \mathfrak{F}}$ is a semi-consequence relation over IPC .

Below we will provide a number of motivating examples for our definitions.

3.4 Remark: It is instructive to compare preservativity with conservativity. Define:

- For $A, B \in \mathfrak{B}$: $A \triangleright^*_{T, X} B \Leftrightarrow \forall C \in X (B \vdash_T C \Rightarrow A \vdash_T C)$,

So $A \triangleright^*_{T, X} B$ means that $T+B$ is conservative over $T+A$ w.r.t. X . For *classical* theories T we have:

$$\star \quad A \triangleright_{T, X} B \Leftrightarrow \neg B \triangleright^*_{T, \neg X} \neg A, \text{ where } \neg X := \{\neg C \mid C \in X\}.$$

Thus classically preservativity is a superfluous notion. Constructively, however, the reduction given in \star does not work.

Note that conservativity as a relation between sentences over a theory is a natural extension of the notion of conservativity between theories. There is no analogue of this for preservativity. \circ

A formula A of \mathfrak{B} is *T-prime* if for any finite set of \mathfrak{B} -formulas Z :

$$A \vdash_T Z \Rightarrow \exists B \in Z A \vdash_T B.$$

Note that if A is T-prime, then $A \not\vdash_T \perp$. A class of formulas X is *weakly T-disjunctive* if every $A \in X$ is equivalent to the disjunction of a finite set of T-prime formulas Y , with $Y \subseteq X$.

For any class of formulas X , let $\text{DISJ}(X)$ be the closure of X under arbitrary conjunctions.

In the next fact we collect a number of noteworthy small facts on preservativity.

3.5 Fact

- i) $\triangleright_{T,X}$ is a semi-consequence relation over T . $\triangleright_{T,X,\mathfrak{F}}$ is a semi-consequence relation over IPC.
- ii) Suppose X is weakly T-disjunctive, then $\triangleright_{T,X}$ and $\triangleright_{T,X,\mathfrak{F}}$ are nearly-consequence relations.
- iii) Suppose X is closed under conjunction, then:

$$C \in X \text{ and } A \triangleright_{T,X} B \Rightarrow (C \rightarrow A) \triangleright_{T,X} (C \rightarrow B).$$
- iv) Let $\text{range}(\mathfrak{F})$ be the union of the ranges of the elements of \mathfrak{F} . Suppose $\text{range}(\mathfrak{F}) \subseteq X$ and X is closed under conjunction, then $\triangleright_{T,X,\mathfrak{F}}$ satisfies B2.
- v) Suppose $A \in X$, then: $A \triangleright_{T,X} B \Leftrightarrow A \vdash_T B$.
- vi) Suppose $X \subseteq Y$ and $\mathfrak{F} \subseteq \mathfrak{G}$, then $\triangleright_{T,Y} \subseteq \triangleright_{T,X}$ and $\triangleright_{T,Y,\mathfrak{G}} \subseteq \triangleright_{T,X,\mathfrak{F}}$.
- vii) $\triangleright_{T,\text{DISJ}(X)} = \triangleright_{T,X}$ and hence $\triangleright_{T,\text{DISJ}(X),\mathfrak{F}} = \triangleright_{T,X,\mathfrak{F}}$.
- viii) Let $\text{id} : \mathfrak{B} \rightarrow \mathfrak{L}(\mathfrak{B})$ be the function with $\text{id}(p) = p$. Then: $\triangleright_{\text{IPC},X,\text{id}} = \triangleright_{\text{IPC},X}$.

Proof: We treat (ii) and (iii). (ii) Suppose X is weakly T-disjunctive and $A \triangleright_{T,X} C$ and $B \triangleright_{T,X} C$. Let E be any element of X . Suppose Y is a finite set of T-prime formulas from X such that E is equivalent to the disjunction of Y . We have:

$$\begin{aligned} E \vdash_T A \vee B &\Rightarrow \bigvee Y \vdash_T A \vee B \\ &\Rightarrow \forall F \in Y F \vdash_T A \vee B \\ &\Rightarrow \forall F \in Y F \vdash_T A \text{ or } F \vdash_T B \\ &\Rightarrow \forall F \in Y F \vdash_T C \\ &\Rightarrow \bigvee Y \vdash_T C \\ &\Rightarrow E \vdash_T C \end{aligned}$$

B1 for $\triangleright_{T,X,\mathfrak{F}}$ is immediate from B1 for $\triangleright_{T,X}$.

iii) Suppose X is closed under conjunction. Suppose $A \triangleright B$. Let $C, E \in X$ and suppose $E \vdash (C \rightarrow A)$. Then $(E \wedge C) \vdash A$ and hence $(E \wedge C) \vdash B$. Ergo $E \vdash (C \rightarrow B)$. \square

3.6 Example 1: Take $X := \{ \top \}$. We have:

For $A, B \in \mathfrak{B}$: $A \triangleright_{T, \top} B \Leftrightarrow (\vdash_T A \Rightarrow \vdash_T B)$,

For $A, B \in \mathcal{L}(\mathfrak{B})$: $A \triangleright_{T, \top, \mathfrak{F}} B \Leftrightarrow \forall f \in \mathfrak{F} (\vdash_T A[f] \Rightarrow \vdash_T B[f])$.

Thus $\triangleright_{T, \top}$ is the relation of *deductive consequence* for T and $\triangleright_{T, \top, \mathfrak{F}}$ is the relation of being a (propositional) *derived rule* for T w.r.t. \mathfrak{F} . We mention the important result due to Rybakov (see Rybakov[92]) that $\triangleright_{IPC, \top, \mathcal{L}(\mathfrak{B})}$ is decidable.

Note that if T has the disjunction property, then $\triangleright_{T, \top}$ and $\triangleright_{T, \top, \mathfrak{F}}$ satisfy B1.

We repeat a well known observation.

3.6.1 Theorem: Suppose for all $C \in \mathfrak{B}$: $\vdash C \Leftrightarrow \forall f \in \mathfrak{F} \vdash_T C[f]$. Moreover suppose $\mathcal{G} \subseteq \mathcal{L}(\mathfrak{B})^{\mathfrak{B}}$ and for any $g \in \mathcal{G}$ and $f \in \mathfrak{F}$: $g \circ f \in \mathfrak{F}$. (We read composition in the order of application, so that $A[g][f] = A[g \circ f]$.) Then: $\triangleright_{T, \top, \mathfrak{F}} \subseteq \triangleright_{IPC, \top, \mathcal{G}}$.

Proof: Suppose $A \triangleright_{T, \top, \mathfrak{F}} B$ and $\vdash A[g]$. Consider any $f \in \mathfrak{F}$. Clearly $\vdash_T A[g][f]$, i.e. $\vdash_T A[g \circ f]$. Since $g \circ f \in \mathfrak{F}$, we find: $\vdash_T B[g \circ f]$. Hence for all $f \in \mathfrak{F}$: $\vdash_T B[g][f]$ and hence: $\vdash B[g]$. \square

Consider e.g. HA. Let \mathfrak{A} be the language of HA, let Σ be the set of Σ -sentences and let $B(\Sigma)$ be the set of Boolean (Brouwerian) combinations of Σ -sentences. We have:

* for all $C \in \mathfrak{A}$: $\vdash C \Leftrightarrow \forall f \in \Sigma^{\mathfrak{B}} \vdash_{HA} C[f]$.

* is De Jongh's completeness theorem for IPC w.r.t. Σ -substitutions in HA. There are many different proofs of *, see e.g. Smoryński [73] or Visser[85] or DV[94]. Let $\mathfrak{F} := B(\Sigma)^{\mathfrak{B}}$ and $\mathcal{G} := \mathcal{L}(\mathfrak{B})^{\mathfrak{B}}$ and $T := HA$ in 3.6.1. We find: $\triangleright_{HA, \top, B(\Sigma)} \subseteq \triangleright_{IPC, \top, \mathcal{L}(\mathfrak{B})}$. In DV[94] (theorem 6.2) it is shown that $\triangleright_{HA, \top, B(\Sigma)} = \triangleright_{IPC, \top, \mathcal{L}(\mathfrak{B})}$. So we have:

✧ $\triangleright_{HA, \top, \mathfrak{A}} \subseteq \triangleright_{HA, \top, B(\Sigma)} = \triangleright_{IPC, \top, \mathcal{L}(\mathfrak{B})} \subseteq \triangleright_{HA, \top, \Sigma}$.

The main result of the present paper implies that $(\neg\neg p \rightarrow p) \triangleright_{HA, \top, \Sigma} (p \vee \neg p)$. On the other hand we have: $IPC \vdash (\neg\neg p \rightarrow p) [p := \neg q]$, but not $IPC \vdash (p \vee \neg p) [p := \neg q]$. So the second inclusion of ✧ is strict. It is open whether the first inclusion of ✧ is strict.

In 3.8 we will see that $\triangleright_{IPC, \top, \mathcal{L}(\mathfrak{B})}$ is extensionally equal to a pure preservativity relation. We will see in 9.2, that $\triangleright_{HA, \top, \Sigma}$ can be alternatively described as a pure preservativity relation.

3.6.2 Reformulation of our main result: The result we are aiming at in this paper is (a)+(e) of the introduction. In our new notation this becomes:

there is an $A^* \in \text{NNIL}$ such that:

- b) $A \triangleright_{\text{HA}, \top, \Sigma} A^*$.
- f) for all $B \in \mathcal{L}(\mathfrak{B})$: $A \triangleright_{\text{HA}, \top, \Sigma} B \Rightarrow \vdash A^* \rightarrow B$.

We show that (b)+(f) is equivalent to:

- g) for all $B \in \mathcal{L}(\mathfrak{B})$: $A \triangleright_{\text{HA}, \top, \Sigma} B \Leftrightarrow \vdash A^* \rightarrow B$.

Thus our main result takes the form of a characterization of $\triangleright_{\text{HA}, \top, \Sigma}$.

Proof: “(b)+(f) \Rightarrow (g)” Suppose (b) and (f). We prove (g). From left to right is trivial. Suppose $\vdash A^* \rightarrow B$. Then by A1: $A^* \triangleright_{\text{HA}, \top, \Sigma} B$. Hence by (b) and A2: $A \triangleright_{\text{HA}, \top, \Sigma} B$.

“(g) \Rightarrow (b)+(f)” is trivial.

$\square((b)+(f) \Leftrightarrow (g))$

By 3.2(i) and (ii) we may conclude that $(.)^*$ gives unique values modulo IPC-provable equivalence and that $(.)^*$ is a co-closure operation on the Heyting algebra of IPC. We will see later that $\triangleright_{\text{HA}, \top, \Sigma}$ satisfies B1 and that, hence, $(.)^*$ commutes with disjunction according to 3.2(iii) \bigcirc

3.7 Example 2: Take $X := \mathfrak{B}$, $\triangleright_{\top, X}$ is *T-provable consequence* or \vdash_{\top} .

3.8 Example 3: We show that the notion of *derived rule for IPC* coincides with a pure preservativity relation.

An $\mathcal{L}(\mathbf{p})$ -formula A is *IPC, \mathbf{p} -exact* if there is an $\mathfrak{f} \in \mathcal{L}(\mathfrak{B})^{\mathbf{p}}$ such that:

- \star for all $B \in \mathcal{L}(\mathbf{p})$: $\vdash B[\mathfrak{f}] \Leftrightarrow A \vdash B$.

Let's say that the class of IPC, \mathbf{p} -exact formulas is $\text{EX}(\mathbf{p})$. Let EX be the union of the $\text{EX}(\mathbf{p})$'s for all finite subsets \mathbf{p} of \mathfrak{B} .

De Jongh & Visser show (see: DV[94], Theorem 2.3) that for any $\mathfrak{f} \in \mathcal{L}(\mathfrak{B})^{\mathbf{p}}$ we can find an $A_{\mathfrak{f}} \in \mathcal{L}(\mathbf{p})$ satisfying \star . It is easy to see that $A_{\mathfrak{f}}$ is unique modulo IPC-provable equivalence. Let \mathbf{q} be the set of variables occurring in the range of \mathfrak{f} . Let \mathfrak{f}^+ be the result of putting $\mathfrak{f}^+(p) := \mathfrak{f}(p)$ if $p \in \mathbf{p}$, $\mathfrak{f}^+(p) := p$ otherwise. Inspection of the argument by de Jongh and Visser reveals that A satisfies the stronger:

- \odot for all $B \in \mathcal{L}(\mathbf{p} \cup (\mathfrak{B}/\mathbf{q}))$: $\vdash B[\mathfrak{f}^+] \Leftrightarrow A \vdash B$.

To see that the variable condition is needed, let e.g. $\mathbf{p} := \{p\}$, $\mathfrak{f} := [p := \neg q]$. It is easy to see that $A_{\mathfrak{f}} = (\neg p \rightarrow p)$. Let $B := (\neg q \rightarrow p)$. $B[\mathfrak{f}^+] = (\neg q \rightarrow \neg q)$ and so $\vdash B[\mathfrak{f}^+]$. On the other hand $(\neg p \rightarrow p) \not\vdash (\neg q \rightarrow p)$. (The proof in DV[94], employs Pitts's Uniform Interpolati-

on Theorem for IPC (see Pitts[92]). Inspecting the proof it can be seen that the variable condition is analogous to the one of \exists -elimination.)

We show: $\triangleright_{\text{IPC}, \top, \mathcal{L}(\mathfrak{P})} = \triangleright_{\text{IPC}, \text{EX}}$.

Suppose $B \triangleright_{\text{IPC}, \top, \mathcal{L}(\mathfrak{P})} C$ and $E \in \text{EX}$, say $E \in \text{EX}(\mathbf{p})$ and let $\mathbf{f} \in \mathcal{L}(\mathfrak{P})^{\mathbf{p}}$ witness the fact that $E \in \text{EX}(\mathbf{p})$. As is easily seen (by permuting \mathfrak{P}) we can arrange that the variables in the elements of the range of \mathbf{f} are disjoint from $\text{PV}(\mathbf{B}) \cup \text{PV}(\mathbf{C})$. Construct \mathbf{f}^+ as above. We find using \odot :

$$E \vdash B \Rightarrow \vdash B[\mathbf{f}^+] \Rightarrow \vdash C[\mathbf{f}^+] \Rightarrow E \vdash C.$$

Conversely suppose $B \triangleright_{\text{IPC}, \text{EX}} C$ and consider $\mathbf{f} \in \mathcal{L}(\mathfrak{P})^{\mathfrak{P}}$. Take $\mathbf{p} := \text{PV}(\mathbf{B}) \cup \text{PV}(\mathbf{C})$ and $\mathbf{g} := \mathbf{f}|_{\mathbf{p}}$ and $\mathbf{A} := \mathbf{A}_{\mathbf{g}}$. Then:

$$\vdash B[\mathbf{f}] \Rightarrow \vdash B[\mathbf{g}] \Rightarrow \mathbf{A} \vdash B \Rightarrow \mathbf{A} \vdash C \Rightarrow \vdash C[\mathbf{g}] \Rightarrow \vdash C[\mathbf{f}].$$

In the cases that $\mathfrak{P} = \emptyset$ and $\mathfrak{P} = \{\mathbf{p}\}$, $\triangleright_{\text{IPC}, \text{EX}}$ is completely understood. The case that $\mathfrak{P} = \{\mathbf{p}\}$ is the subject of De Jongh[82]. He shows that the exact formulas in one variable are precisely: \mathbf{p} , $\neg \mathbf{p}$, $\neg \neg \mathbf{p}$, $\neg \neg \mathbf{p} \rightarrow \mathbf{p}$, \top . It follows that $\text{DISJ}(\text{EX})$ is finite (mod IPC). If we set $\Phi(\mathbf{A}) :=$ the disjunction of all EX-formulas that IPC-imply \mathbf{A} , we get: $\mathbf{A} \triangleright_{\text{IPC}, \text{EX}} \mathbf{B} \Leftrightarrow \Phi(\mathbf{A}) \vdash \mathbf{B}$. \circ

3.9 Remark: It is not difficult to see that every $\text{NNIL}(\mathbf{p})$ -formula can be written as a finite disjunction of prime $\text{NNIL}(\mathbf{p})$ -formulas. In DV[94] (Theorem 2.8) it is shown that prime $\text{NNIL}(\mathbf{p})$ -formulas are IPC, \mathbf{p} -exact. So modulo IPC-provable equivalence: $\text{NNIL}(\mathbf{p}) \subseteq \text{DISJ}(\text{EX}(\mathbf{p}))$. Ergo by 3.5(vii) and 3.8:

$$\triangleright_{\text{IPC}, \top, \mathcal{L}(\mathfrak{P})} = \triangleright_{\text{IPC}, \text{EX}} = \triangleright_{\text{IPC}, \text{DISJ}(\text{EX})} \subseteq \triangleright_{\text{IPC}, \text{NNIL}}.$$

We leave it to the reader to show that the last inclusion is proper. (We will show in 5.2 and 9.2 that $\triangleright_{\text{IPC}, \text{NNIL}} = \triangleright_{\text{HA}, \top, \Sigma}$. So our result here coheres with \clubsuit of 3.6.) \circ

3.10 Example 4: Suppose $\mathcal{Q} \subseteq \mathfrak{P}$. Consider $\triangleright_{\text{IPC}, \mathcal{L}(\mathcal{Q})}$. By 2.5, $\mathcal{L}(\mathcal{Q})$ is weakly IPC-disjunctive. Hence, by 3.5(ii), $\triangleright_{\text{IPC}, \mathcal{L}(\mathcal{Q})}$ is a nearly consequence relation. Note that $\triangleright_{\text{IPC}, \mathcal{L}(\mathcal{Q})}$ also satisfies:

$$\text{if } \mathbf{C} \in \mathcal{L}(\mathcal{Q}): \mathbf{A} \triangleright \mathbf{B} \Rightarrow (\mathbf{C} \rightarrow \mathbf{A}) \triangleright (\mathbf{C} \rightarrow \mathbf{B}),$$

$$\text{if } \mathbf{p} \notin \mathcal{L}(\mathcal{Q}): (\mathbf{p} \rightarrow \mathbf{A}) \triangleright \mathbf{A}[\mathbf{p} := \top].$$

Andrew Pitts, in his Pitts[92], shows that for every $\mathbf{A} \in \mathcal{L}(\mathfrak{P})$, there is a formula $(\forall \mathbf{A}) \in \mathcal{L}(\mathcal{Q})$, such that for all $\mathbf{B} \in \mathcal{L}(\mathfrak{P})$: $\mathbf{A} \triangleright_{\text{IPC}, \mathcal{L}(\mathcal{Q})} \mathbf{B} \Leftrightarrow \forall \mathbf{A} \vdash \mathbf{B}$. We can read $\forall \mathbf{A}$ informally as the result of universally quantifying out the propositional variables not in \mathcal{Q} . Note that we may apply 3.2 to \forall .

3.11 List of results on preservativity: We close this section by listing a number of relationships between various preservativity relations.

- 1) $\triangleright_{\text{IPC}, \mathfrak{L}(\mathfrak{P})}$ is the minimal preservativity relation and:
 $\triangleright_{\text{IPC}, \mathfrak{L}(\mathfrak{P})} \subset^{(\alpha)} \triangleright_{\text{HA}, \top, \mathfrak{U}} \subseteq^{(\beta)} \triangleright_{\text{IPC}, \text{EX}} \subset^{(\gamma)} \triangleright_{\text{IPC}, \text{NNIL}} \subset \triangleright_{\text{IPC}, \top}$.
 - α) Non-identity follows from the presence of non-trivial derived rules for HA, like IP, i.e. $(\neg p \rightarrow (q \vee r)) \triangleright_{\text{HA}, \top, \mathfrak{U}} ((\neg p \rightarrow q) \vee (\neg p \rightarrow r))$.
 - β) Immediate from (3) below.
 - γ) Non-identity follows by: $(\neg \neg p \rightarrow p) \triangleright_{\text{IPC}, \text{NNIL}} (p \vee \neg p)$ and the fact that $(\neg \neg p \rightarrow p)$ is IPC, p -exact by the substitution $[p := \neg q]$. Hence we do not have $(\neg \neg p \rightarrow p) \triangleright_{\text{IPC}, \text{EX}} (p \vee \neg p)$.
- 2) $\triangleright_{\text{IPC}, \mathfrak{L}(\mathfrak{P}), \mathfrak{L}(\mathfrak{P})} = \triangleright_{\text{IPC}, \mathfrak{L}(\mathfrak{P})} =^{(\delta)} \triangleright_{\text{HA}, \mathfrak{U}, \mathfrak{U}} =^{(\delta)} \triangleright_{\text{HA}, \mathfrak{U}, \Sigma} =^{(\delta)} \triangleright_{\text{HA}^*, \mathfrak{U}, \mathfrak{U}} =^{(\delta)} \triangleright_{\text{HA}^*, \mathfrak{U}, \Sigma} =^{(\epsilon)} \triangleright_{\text{HA}^*, \top, \mathfrak{U}} =^{(\epsilon)} \triangleright_{\text{HA}^*, \top, \Sigma}$.
 - δ) HA* is introduced in section 6. All the identities follow immediately from De Jongh's Completeness Theorem for Σ -substitutions for HA and HA*. See e.g DV[94].
 - ϵ) The identity of $\triangleright_{\text{HA}^*, \top, \mathfrak{U}}$ and $\triangleright_{\text{HA}^*, \top, \Sigma}$ with $\triangleright_{\text{IPC}, \mathfrak{L}(\mathfrak{P})}$, is corollary 5.7 of DV[94]
- 3) $\triangleright_{\text{IPC}, \text{EX}} =^{(\eta)} \triangleright_{\text{IPC}, \top, \mathfrak{L}(\mathfrak{P})} =^{(\zeta)} \triangleright_{\text{HA}, \top, \text{B}(\Sigma)}$.
 - η) See 3.8.
 - ζ) This is theorem 6.2 of DV[94].
- 4) $\text{NNIL} =^{(\vartheta)} \text{ROB} =^{(\vartheta)} \text{f-ROB (mod IPC)}$, and
 $\triangleright_{\text{IPC}, \text{NNIL}} =^{(\iota)} \triangleright_{\text{HA}, \Sigma, \Sigma} =^{(\iota)} \triangleright_{\text{HA}, \top, \Sigma}$.
 - ϑ) ROB and f-ROB are defined in section 4 below. The result is proved in Visser[85] and by a similar proof in section 5 of this paper. A different proof, due to Johan van Benthem, is contained in Van Benthem[91]. A version of van Benthem's proof is contained in DRVV[94].
 - ι) See section 9.

In the next section we study robust formulas, as a preliminary to the specification of the NNIL-algorithm in section 5.

4 Robust formulas: Consider a \mathfrak{P} -models \mathbb{K} and \mathbb{M} . We say that:

- $\mathbb{K} \subseteq \mathbb{M} :\Leftrightarrow \exists f: \mathbb{K} \rightarrow \mathbb{M}, f \text{ is injective and}$
 $\leq_{\mathbb{K}} \circ f \subseteq f \circ \leq_{\mathbb{M}} \text{ and } \forall k \in \mathbb{K}, p \in \mathfrak{P} (k \models_{\mathbb{K}} p \Leftrightarrow f(k) \models_{\mathbb{M}} p)$.

Note that, modulo the identification of the elements of \mathbb{K} with their f -images in \mathbb{M} , ' $\mathbb{K} \subseteq \mathbb{M}$ ' means that \mathbb{K} is a submodel of \mathbb{M} .

- $A \in \text{ROB} :\Leftrightarrow A \text{ is robust} :\Leftrightarrow \forall \mathbb{M} (\mathbb{M} \models A \Rightarrow \forall \mathbb{K} \subseteq \mathbb{M} \mathbb{K} \models A)$

We will let $\sigma, \sigma', \tau, \dots$ range over ROB. It is easy to see that $\text{NNIL} \subseteq \text{ROB}$. In section 5 we will see that modulo IPC-provable equivalence each robust formula is in NNIL.

In this section we will prove that $\triangleright_{\text{IPC}, \text{ROB}}$ is a σ -relation. This result will be our main tool (in section 5) for proving the NNIL-algorithm to be correct. Let's take as a local convention of this section: $\triangleright := \triangleright_{\text{IPC}, \text{ROB}}$.

Clearly we have: $\mathfrak{R} \subseteq \text{ROB}$ and ROB is closed under conjunction. So it follows that \triangleright satisfies A1, A2, A3, B2. To verify B1, by 3.5(ii), the following theorem suffices:

4.1 Theorem: ROB is weakly disjunctive.

Proof: Consider any $\sigma \in \text{ROB}$. We write σ in disjunctive normal form D_σ as in 2.5.1. Consider any disjunct $C(\Delta)$ of D_σ : here, as in 2.5.1, Δ is a \subseteq -minimal $\text{SUB}(\sigma)$ -saturated set with $\sigma \in \Delta$ and $C(\Delta)$ is the conjunction of the atoms and implications in Δ . We claim that $C(\Delta)$ is robust. Consider any models $\mathbb{K} \subseteq \mathbb{M} \models C(\Delta)$. Trivially: $\mathbb{M} \models \Delta$. By the Push Down Lemma 2.3:

$$\text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}[\Delta], \mathbb{M}) \models \Delta.$$

Hence:

$$\text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}[\Delta], \mathbb{M}) \models \sigma.$$

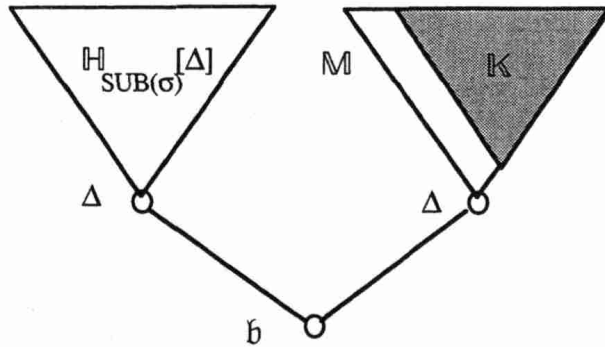
Now clearly:

$$\text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}(\Delta), \mathbb{K}) \subseteq \text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}(\Delta), \mathbb{M}).$$

By the stability of σ we get: $\text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}(\Delta), \mathbb{K}) \models \sigma$. Consider

$$\Gamma := \{G \in \text{SUB}(\sigma) \mid \text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}(\Delta), \mathbb{K}) \models G\}.$$

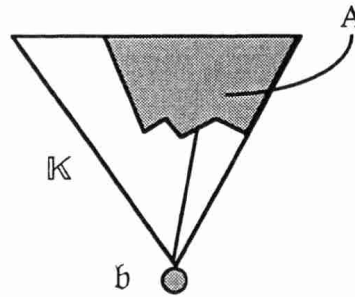
Clearly $\Gamma \subseteq \Delta$ and Γ is $\text{SUB}(\sigma)$ -saturated and $\sigma \in \Gamma$. By the \subseteq -minimality of Δ we find: $\Gamma = \Delta$. Hence $\text{Glue}(\mathbb{H}_{\text{SUB}(\sigma)}(\Delta), \mathbb{K}) \models C(\Delta)$ and so $\mathbb{K} \models C(\Delta)$. Ergo $C(\Delta)$ is robust. \square



4.2 Theorem: \triangleright is closed under B3.

Proof: Let X be a finite set of implications and let $Y := \{C \mid (C \rightarrow D) \in X\} \cup \{B\}$. Let $A := \bigwedge X$. We have to show: $(A \rightarrow B) \triangleright [A]Y$. The proof is by contraposition. Consider any $\sigma \in \text{ROB}$ and suppose: $\sigma \not\models [A]Y$. Let $\mathbb{K} = \langle K, b, \leq, \models, \models' \rangle$ be a rooted model such that $\mathbb{K} \models \sigma$ and $\mathbb{K} \not\models \bigvee [A]Y$, i.e. for all $E \in Y$: $\mathbb{K} \not\models [A]E$.

Let $\mathbb{K}' := \mathbb{K}\{A\} := \langle K', b, \leq', \models', \models' \rangle$, be the full submodel of \mathbb{K} on $K' := \{b\} \cup \{k' \in K \mid k' \models A\}$. (A submodel is *full* if the new ordering relation is the restriction of the old ordering relation to the new worlds.) Note that $\{k' \in K \mid k' \models A\}$ is upwards closed and that on $\{k' \in K \mid k' \models A\}$ the old and the new forcing coincide. Moreover on this class $[A]G$ is equivalent with G .



4.2.1 Claim: For all F : $b \models' F \Rightarrow b \models [A]F$.

Proof of the claim: The proof is by a simple induction on F . The cases of atoms disjunction and conjunction are trivial. Suppose F is an implication and $b \models' F$. Then certainly: for all $k \in K$ ($k \models A \Rightarrow k \models' F$). Since on the k with $k \models A$, \models and \models' coincide, we find: $b \models (A \rightarrow F)$, i.e. $b \models [A]F$. \square (Claim)

We return to the main proof. Remember that:

$b \models \sigma$ and $b \not\models [A]C$ for all C with $(C \rightarrow D) \in X$ and $b \not\models [A]B$.

We show that $b \models' \sigma$ and $b \models' A$ and $b \not\models' B$.

It is immediate that $b \models' \sigma$, since \mathbb{K}' is a submodel of \mathbb{K} and σ is robust.

Remember that A is the conjunction of the $(C \rightarrow D)$ in X . So it is sufficient to show that for each $(C \rightarrow D)$ in X : $b \models' (C \rightarrow D)$. Consider any $(C \rightarrow D) \in X$ and any $k' \geq' b$ with $k' \models' C$. Since $b \not\models [A]C$, we have by the claim: $b \not\models' C$. So $k' \neq b$. But then $k' \models A$, hence $k' \models' A$ and thus $k' \models' (C \rightarrow D)$. We may conclude: $k' \models' D$ and hence $b \models' (C \rightarrow D)$.

From $b \not\models [A]B$ and the claim we have immediately that: $b \not\models' B$. \square

All the proofs in this section also work when we replace ROB by f-ROB, the set of formulas preserved by full submodels. Note that $\text{ROB} \subseteq \text{f-ROB}$. Our result in section 5 will imply: $\text{ROB} = \text{f-ROB} = \text{NNIL}$ (modulo IPC-provable equivalence).

5 The NNIL-algorithm: In this section we produce the algorithm that is the main tool of this paper. The existence of the algorithm establishes the following theorem.

5.1 Theorem: For all A there is an $A^* \in \text{NNIL}(\text{PV}(A))$ such that for all σ -relations \triangleright : $A \triangleright A^*$ and $A^* \vdash A$.

For convenience we reproduce the defining properties of σ -relations here.

- A1 $A \vdash B \Rightarrow A \triangleright B$
- A2 $A \triangleright B \text{ and } B \triangleright C \Rightarrow A \triangleright C$
- A3 $C \triangleright A \text{ and } C \triangleright B \Rightarrow C \triangleright (A \wedge B)$
- B1 $A \triangleright C \text{ and } B \triangleright C \Rightarrow (A \vee B) \triangleright C$
- B2 $A \triangleright B \Rightarrow (p \rightarrow A) \triangleright (p \rightarrow B)$
- B3 Let X be a finite set of implications and let $Y := \{C \mid (C \rightarrow D) \in X\} \cup \{B\}$. Take $A := \bigwedge X$. Then: $(A \rightarrow B) \triangleright [A]Y$

Before proceeding with the proof of 5.1 we interpolate a corollary. Let's write "mod IPC" for "modulo IPC-provable equivalence".

5.2 Corollary

- i) Let A and A^* be as promised by 5.1. Then we have:

$$A \triangleright_{\text{ROB}} B \Leftrightarrow A^* \vdash B.$$
It follows that $(.)^*$ has the properties promised in 3.2.
- ii) $\text{NNIL} = \text{ROB} = \text{f-ROB}$ (mod IPC)
- iii) $\triangleright_{\text{ROB}}$ is the minimal σ -relation. It follows that A1-A3, B1-B3 axiomatizes $\triangleright_{\text{ROB}}$.

Proof: (i) Suppose $A \triangleright_{\text{ROB}} B$. Since $A^* \vdash A$, it follows by A1 that $A^* \triangleright_{\text{ROB}} A$ and, hence, by A2 that $A^* \triangleright B$. Since $A^* \in \text{NNIL} \subseteq \text{ROB}$, we have, by 3.5(v), that $A^* \vdash B$. Conversely if $A^* \vdash B$, then $A^* \triangleright_{\text{ROB}} B$. Since $\triangleright_{\text{ROB}}$ is a σ -relation, we have: $A \triangleright_{\text{ROB}} A^*$ and hence by A2: $A \triangleright_{\text{ROB}} B$.

(ii) Consider $A \in \text{f-ROB}$. We have $A \triangleright_{\text{f-ROB}} A^*$, so by 3.5(v): $A \vdash A^*$. Since also $A^* \vdash A$, we find: $\vdash A \leftrightarrow A^*$. Ergo $\text{f-ROB} \subseteq \text{NNIL}$ (mod IPC). Since obviously $\text{NNIL} \subseteq \text{ROB} \subseteq \text{f-ROB}$, we are done.

(iii) Let \triangleright be any σ -relation. We have by (i), A1, 5.1, A2:

$$A \triangleright_{\text{ROB}} B \Rightarrow A^* \vdash B \Rightarrow A^* \triangleright B \Rightarrow A \triangleright B. \quad \square$$

5.2(ii) was proved by purely model theoretical means by Johan van Benthem. See his Van Benthem[91] or DRVV[94]. The advantage of van Benthem's proof is its relative simplicity and the fact that the method employed easily generalizes. The advantage of the present method is the extra information it produces, like 5.2(iii) and its usefulness in the arithmetical case, see sections 8, 9 and 10. It is an open question whether van Benthem's proof can be adapted to the arithmetical case.

Proof of 5.1: We introduce an ordinal measure \circ of complexity on formulas as follows:

- $I(D) := \{E \in \text{SUB}(D) \mid E \text{ is an implication}\},$
- $i(D) := \max \{|I(E)| \mid E \in I(D)\},$
- $c(D) := \text{the number of occurrences of logical connectives in } D,$
- $\circ(D) := \omega \cdot i(D) + c(D).$

Note that we count *occurrences* of connectives for c and *types* of implications, not occurrences, for i !

We say that an occurrence of E in D is an *outer occurrence* if this occurrence is not in the scope of an implication.

We prove 5.1 by induction on \circ . Consider any σ -relation \triangleright . We will exhibit the properties of \triangleright that are used between square brackets.

α **Atoms:** [A1] Suppose A is an atom. Take $A^* := A$.

β **Conjunction:** [A1,A2,A3]: Suppose $A = (B \wedge C)$. Clearly $\circ(B) < \circ(A)$ and $\circ(C) < \circ(A)$. Take $A^* := B^* \wedge C^*$. Clearly $A^* \in \text{NNIL}(\text{PV}(A))$. The verification of the desired properties is trivial.

γ **Disjunction:** [A1,A2,B1] Suppose $A = (B \vee C)$. Clearly $\circ(B) < \circ(A)$ and $\circ(C) < \circ(A)$. Take $A^* := B^* \vee C^*$. Clearly $A^* \in \text{NNIL}$. The verification of the desired properties is trivial.

δ **Implication:** [A1,A2,A3,B2,B3] Suppose $A = (B \rightarrow C)$. We split into several cases.

δ1 Outer conjunction in the consequent: $[A1, A2, A3]$ Suppose C has an outer occurrence of a formula $D \wedge E$. Pick any $J(q)$ be such that:

- q is a fresh variable,
- q occurs precisely once in J ,
- q is not in the scope of an implication in J ,
- $C = J[q := D \wedge E]$.

Let $C_1 := J[q := D]$, $C_2 := J[q := E]$. As is easily seen C is IPC-provably equivalent to $C_1 \wedge C_2$. Let $A_1 := (B \rightarrow C_1)$ and $A_2 := (B \rightarrow C_2)$. Clearly A is IPC-provably equivalent to $A_1 \wedge A_2$. We prove: $\circ(A_i) < A$ for $i=0,1$. Since it is clear that $c(A_i) < c(A)$, it is sufficient to show that $i(A_i) \leq i(A)$. Since A and the A_i are implications we have to show that: $II(A_i) \leq II(A)$. We treat the case that $i=1$. It is sufficient to construct an injective mapping from $I(A_1)$ to $I(A)$. Consider any implication F in $I(A_1)$. If $F = A_1$, we map F to A . Otherwise $F \in I(B)$ or $F \in I(J)$ or $F \in I(D)$ (since q does not occur in the scope of an implication). In all three cases we can map F to itself. Since A_1 cannot be in $I(B)$ or $I(J)$ or $I(D)$, our mapping is injective. The case that $i=2$ is similar. Set $A^* := (A_1^* \wedge A_2^*)$.

δ2 Outer disjunction in the antecedent: $[A1, A2, A3]$ This case is completely analogous to the previous one.

If A has no outer disjunction in the antecedent and no outer conjunction in the consequent, then B is a conjunction of atoms and implications and C is a disjunction of atoms and implications. It is easy to see, that applications of associativity, commutativity and idempotency to the conjunction in the antecedent or to the disjunction in the consequent do not raise \circ . So we can safely write: $B = \bigwedge X$ and $C = \bigvee Y$, where X and Y are finite sets of atoms or implications. This leads us to the following case.

δ3 B is a conjunction of atoms and implications, C is a disjunction of atoms and implications $[A1, A2, A3, B2, B3]$

δ3.1 X contains an atom: $[A1, A2, B2]$

δ3.1.1 X contains a propositional variable, say p : $[A1, A2, B2]$ Consider: $D := \bigwedge (X / \{p\})$ and $E := (D \rightarrow C)$. Clearly $\vdash A \leftrightarrow (p \rightarrow E)$ and $\circ(E) < \circ(A)$. Put $A^* := (p \rightarrow E^*)$. Evidently $(p \rightarrow E^*)$ is in $NNIL(PV(A))$. We have $E^* \vdash E$ by the Induction Hypothesis and hence: $(p \rightarrow E^*) \vdash (p \rightarrow E) \vdash A$. We have $E \triangleright E^*$ by the Induction Hypothesis. It follows by B2 that: $(p \rightarrow E) \triangleright (p \rightarrow E^*)$. From $A \vdash (p \rightarrow E)$, we have by A1: $A \triangleright (p \rightarrow E)$. Hence by A2: $A \triangleright (p \rightarrow E^*)$.

§3.1.2 X contains \top : [A1,A2] Left to the reader.

§3.1.3 X contains \perp : [A1] Left to the reader.

§3.2 X contains no atoms: [A1,A2,A3,B3] This case -the last one- is the truly difficult one. To motivate it, let's solve the difficulties one by one. Let's first look at an example.

Example 1: Consider $(p \rightarrow q) \rightarrow r$. B3 gives us: $((p \rightarrow q) \rightarrow r) \triangleright (p \vee r)$. However, we do not have: $(p \vee r) \vdash ((p \rightarrow q) \rightarrow r)$. We can repair this by noting that $(p \rightarrow q) \rightarrow r \vdash (q \rightarrow r)$ and $(p \rightarrow q) \wedge (p \vee r) \vdash ((p \rightarrow q) \rightarrow r)$. So the full solution of our example is as follows.

We have: $((p \rightarrow q) \rightarrow r) \triangleright ((q \rightarrow r) \wedge (p \vee r))$, by:

- | | | |
|----|--|--------|
| a) | $((p \rightarrow q) \rightarrow r) \triangleright (p \vee r)$ | B3 |
| b) | $((p \rightarrow q) \rightarrow r) \vdash (q \rightarrow r)$ | IPC |
| c) | $((p \rightarrow q) \rightarrow r) \triangleright (q \rightarrow r)$ | A1 |
| d) | $((p \rightarrow q) \rightarrow r) \triangleright ((q \rightarrow r) \wedge (p \vee r))$ | a,c,A3 |

Moreover: $((q \rightarrow r) \wedge (p \vee r)) \vdash ((p \rightarrow q) \rightarrow r)$ and $((q \rightarrow r) \wedge (p \vee r)) \in \text{NNIL}(\text{PV}((p \rightarrow q) \rightarrow r))$.

○

We implement this idea for the general case. There will be a problem with \circ , but we will postpone its discussion until we run into it. A is of the form $B \rightarrow C$, where B is a conjunction of a finite set of implications X and C is a disjunction of a finite set of atoms or implications Y. For any $D := (E \rightarrow F) \in X$, let:

- $B \downarrow D := \wedge((X/\{D\}) \cup \{F\})$.

Clearly and $\circ((B \downarrow D) \rightarrow C) < \circ(A)$.

Let $Z := \{E \mid (E \rightarrow F) \in X\} \cup \{C\}$. Put: $A_0 := \vee [B]Z$. We will show that our problem reduces to the question whether A_0^* exists. So for the moment we pretend it does. We have:

- | | | |
|-----|--|-----------|
| a) | $A \triangleright A_0$ | B3 |
| a1) | $A_0 \triangleright A_0^*$ | ASS |
| a2) | $A \triangleright A_0^*$ | a,a1,A2 |
| b) | for all $D \in X$: $A \vdash ((B \downarrow D) \rightarrow C)$ | IPC |
| c) | for all $D \in X$: $A \triangleright ((B \downarrow D) \rightarrow C)$ | A1 |
| c1) | for all $D \in X$: $((B \downarrow D) \rightarrow C) \triangleright ((B \downarrow D) \rightarrow C)^*$ | IH |
| c2) | for all $D \in X$: $A \triangleright ((B \downarrow D) \rightarrow C)^*$ | c,c1,A2 |
| d) | $A \triangleright \wedge\{((B \downarrow D) \rightarrow C)^* \mid D \in X\} \wedge A_0^*$ | a2,c2,A3. |

It is clear that $\bigwedge\{((B \downarrow D) \rightarrow C) * ID \in X \wedge A_0^* \in \text{NNIL}(\text{PV}(A))\}$. We show that:

$$\bigwedge\{((B \downarrow D) \rightarrow C) * ID \in X \wedge A_0^* \vdash A.$$

It is sufficient to show:

$$\bigwedge\{((B \downarrow D) \rightarrow C) * ID \in X \wedge [B]E \vdash A \text{ for each } E \in Z.$$

In case $E=C$, we are immediately done by: $[B]C \vdash B \rightarrow C$. Suppose $(E \rightarrow F) \in X$ for some F . Reason in IPC.

Suppose $\bigwedge\{((B \downarrow D) \rightarrow C) * ID \in X\}$, $[B]E$ and B . We want to derive C . We have: $((\bigwedge(X/\{E \rightarrow F\}) \wedge F) \rightarrow C)$, $[B]E$ and B . From B we find $\bigwedge(X/\{E \rightarrow F\})$ and $(E \rightarrow F)$. From B and $[B]E$, we derive E . From E and $(E \rightarrow F)$ we get F . Finally we infer from $((\bigwedge(X/\{E \rightarrow F\}) \wedge F) \rightarrow C)$, $\bigwedge(X/\{E \rightarrow F\})$ and F the desired conclusion C .

So the only thing left to do is to show that A_0^* exists. If $i(A_0)=0$, we are easily done. Suppose $i(A_0)>0$. If for each E in Z with $i(E)>0$, we would have: $i([B]E)<i(A)$, it would follow that $i(A_0)<i(A)$. Hence we would be done by the Induction Hypothesis.

We study some examples. These examples show that we cannot generally hope to get $i([B]E)<i(A)$. The examples will, however, suggest a way around the problem: we produce a logical equivalent, say Q , of $[B]E$, such that $i(Q)<i(A)$. So we replace A_0 by the disjunction R of the Q 's, which is logically equivalent and has $i(R)<i(A_0)$. Put $A_0^*:=R^*$.

Example 2: Consider $(p \rightarrow q) \rightarrow (r \rightarrow s)$. Take $E:=C=(r \rightarrow s)$. We have:

$$[B]E = [p \rightarrow q](r \rightarrow s) = (p \rightarrow q) \rightarrow (r \rightarrow s).$$

So no simplification is obtained. However $[B]E$ is IPC-provably equivalent to: $((r \wedge (p \rightarrow q)) \rightarrow s)$, which has lower i . By A1, A2 we can put $([B]E)^* := ((r \wedge (p \rightarrow q)) \rightarrow s)^*$.

We could have applied the reduction before going to $[B]E$. We did not choose to do so for reasons of uniformity of treatment. \bigcirc

Example 3: Consider $((p \rightarrow q) \rightarrow r) \rightarrow s$. Take $E:=(p \rightarrow q)$. We find:

$$[B]E = [(p \rightarrow q) \rightarrow r](p \rightarrow q) = (((p \rightarrow q) \rightarrow r) \rightarrow (p \rightarrow q)).$$

$[B]E$ is IPC-provably equivalent to $((p \wedge (q \rightarrow r)) \rightarrow q)$, which has lower i . \bigcirc

Consider $[B]E$ for $E \in Z = \{E \mid (E \rightarrow F) \in X\} \cup \{C\}$. $[B]E$ is the result over replacing outer implications $(G \rightarrow H)$ of E by $(B \rightarrow (G \rightarrow H))$. If there is no outer implication in E , we find: $[B]E = E$ and $i([B]E) = 0$. In this case $[B]E$ is implication free and hence a fortiori in NNIL. We can put $([B]E)^* := E$. Suppose E has an outer implication.

Consider any outer implication $(G \rightarrow H)$ of E . We first show: $i(B \rightarrow (G \rightarrow H)) \leq i(B \rightarrow C)$. We define an injection from $I(B \rightarrow (G \rightarrow H))$ to $I(B \rightarrow C)$. Any implication occurring in $I(B)$ or $I(G \rightarrow H)$ can be mapped to itself in $I(B \rightarrow C)$. None of these implications has as image $(B \rightarrow C)$, since $(G \rightarrow H) \in I(B) \cup I(C)$. So we can send $(B \rightarrow (G \rightarrow H))$ to $(B \rightarrow C)$.

The next step is to replace $(B \rightarrow (G \rightarrow H))$ by an IPC-provably equivalent formula K with lower i . Let B' be the result of replacing all occurrences of $(G \rightarrow H)$ in B by H and let: $K := ((G \wedge B') \rightarrow H)$. Clearly K is provably equivalent to $(B \rightarrow (G \rightarrow H))$. We show that $i(K) < i(B \rightarrow (G \rightarrow H))$.

We define a non-surjective injection from $I(K)$ to $I(B \rightarrow (G \rightarrow H))$. Consider an implication M in $I(K)$. If $M=K$ we send it to $(B \rightarrow (G \rightarrow H))$. Suppose $M \neq K$, i.e. $M \in I(B') \cup I(G) \cup I(H)$. A subformula N of $I(B) \cup I(G) \cup I(H)$ is a *predecessor* of M if M is the result of replacing all occurrences of $(G \rightarrow H)$ in N by H . Clearly M has at least one predecessor and any predecessor of M must be an implication. Send M to one of its shortest predecessors. Clearly our function is injective, since two different implications cannot share a predecessor. Finally $(G \rightarrow H)$ cannot be in the range of our injection. If it were, we would have $M=H$, but then H would be a shorter predecessor. A contradiction.

Replace every outer implication in $[B]E$ by an equivalent with lower i . The result is the desired Q . □

5.3 Remarks: The algorithm given with our proof is non-deterministically specified. However by 5.2 the result is unique modulo provable equivalence. I didn't try to make the algorithm optimally quick. In practice the best thing to do is: follow the main line of the algorithm, but, locally, apply ad hoc simplifications.

In Visser[85] a simple adaptation of the NNIL algorithm is given for $\triangleright_{IPC, \top}$. Here the algorithm computes not a value in NNIL, but a value in $\{\top, \perp\}$. Thus we obtain an algorithm that checks for provability. It is easy to use the present algorithm for this purpose too, since there is a p-time algorithm to decide IPC-provability of $NNIL_0$ formulas, i.e. formulas with only propositional variables in front of arrows. The outputs of our algorithm are in $NNIL_0$. It has been shown by Richard Statman (see Statman[79]) that checking whether a formula is IPC-provable is p-space complete. This puts an absolute bound on what an algorithm like ours can do.

5.4 Sample computations (using some obvious short-cuts)

$$\begin{aligned}
& (\neg\neg p \rightarrow p) \rightarrow (p \vee \neg p) \equiv \\
& (p \rightarrow (p \vee \neg p)) \wedge (((\neg\neg p \rightarrow p) \rightarrow \neg p) \vee ((\neg\neg p \rightarrow p) \wedge p) \vee ((\neg\neg p \rightarrow p) \wedge \neg p)) \equiv \\
& \neg\neg p \vee p \vee \neg p \equiv \\
& ((\perp \rightarrow p) \wedge ((\neg p) \vee p \vee (\neg\neg p) \vee \perp)) \vee p \vee \neg p \equiv \\
& p \vee \neg p.
\end{aligned}$$

$$\begin{aligned}
& ((p \rightarrow q) \rightarrow r) \rightarrow s \equiv \\
& (r \rightarrow s) \wedge (([p \rightarrow q] \rightarrow r) \rightarrow s) \vee ([p \rightarrow q] \rightarrow r) \wedge s \equiv \\
& (r \rightarrow s) \wedge ((p \wedge (q \rightarrow r)) \rightarrow r) \vee s \equiv \\
& (r \rightarrow s) \wedge (p \rightarrow ((q \rightarrow r) \rightarrow r)) \vee s \equiv \\
& (r \rightarrow s) \wedge (p \rightarrow ((r \rightarrow q) \wedge ([q \rightarrow r] \vee [q \rightarrow r] \vee q))) \vee s \equiv \\
& (r \rightarrow s) \wedge ((p \rightarrow q) \vee s).
\end{aligned}$$

6 Basic facts and notations in Arithmetic

6.1 Arithmetical Theories: The arithmetical theories T considered in this paper are RE theories in \mathcal{U} , the language of HA. These theories all extend $i\text{-I}\Delta_0 + \text{Exp}$: the constructive theory of Δ_0 -induction with the axiom expressing that the exponentiation function is total. $i\text{-I}\Delta_0 + \text{Exp}$ is finitely axiomatizable; we stipulate that a fixed finite axiomatization is employed. We will use \Box_T for the formalization of provability in T . Suppose A is a formula. $\Box_T A$ means $\text{Prov}_T(t(\mathbf{x}))$, where Prov_T is the arithmetization of the provability predicate of T and where $t(\mathbf{x})$ is the term 'the Gödelnumber of the result of substituting the Gödelnumbers of the numerals of the \mathbf{x} 's for the variables in A '.

We illustrate the above by an example. We suppose Gödelnumbers are assigned as follows:

(11
)	12
=	15
S	8
0	3
+	19

$*$ is the arithmetization of the syntactical operation of concatenation. We use underlining for our external numeral function. num is the arithmetization of the numeral function. We have e.g.: $\text{HA} \vdash \text{num}(\underline{3}) = \underline{8*8*8*3}$. And:

$$\Box_T(x=x) \text{ means: } \text{Prov}_T(\underline{11*\text{num}(x)*15*\text{num}(x)*12}).$$

Our notational convention evidently introduces a scope ambiguity. What is $\Box_T((x+y)=z)$ going to mean?

a) $\text{Prov}_T(\underline{11}*\text{num}(x+y)*\underline{15}*\text{num}(z)*\underline{12})$ (wide scope)

or:

b) $\text{Prov}_T(\underline{11}*\underline{11}*\text{num}(x)*\underline{19}*\text{num}(y)*\underline{12}*\underline{15}*\text{num}(z)*\underline{12})$ (small scope)

Fortunately by standard metamathematical results, we know that as long as the terms we employ stand for T-provably total recursive functions the different readings are provably equivalent. So (a) and (b) are provably equivalent. In this paper we will only employ terms for primitive recursive functions, so the ambiguity is harmless.

We write T_n for the theory axiomatized by the finitely many axioms of $i\text{-IA}_0 + \text{Exp}$, plus the axioms of T, which are smaller than n in the standard Gödelnumbering. We write $\text{Prov}_{T,n}$ for the formalization of provability in T_n . We consider $\text{Prov}_{T,n}$ as a form of *restricted provability* in T. The following well-known fact is quite important:

6.1.1 Fact: Suppose T is an RE extension of HA in the language of HA. Then T is essentially reflexive (verifiably in HA). I.e. we have:

For all n and for all \mathcal{U} -formulas A with free variables \mathbf{x} : $T \vdash \forall \mathbf{x} (\Box_{T,n} A \rightarrow A)$.

And (using UC for: 'the universal closure of') even:

$HA \vdash \forall x \forall A \in \text{FOR}_{\mathcal{U}} \Box_T \text{UC}(\Box_{T,x} A \rightarrow A)$.

Proofsketch: The proof is roughly as follows. Ordinary cut-elimination for constructive predicate logic (or normalization in case we have a natural deduction system) can be formalized in HA. Reason in HA. Let a number x and a formula A be given. Introduce a measure of complexity on arithmetical formulas that counts both the depth of quantifiers and of implications. Find y such that both the axioms of T_x and A have complexity $< y$. We can construct a truthpredicate True_y for formulas of complexity $< y$. We have: $\Box_T \text{UC}(\text{True}_y A \rightarrow A)$. Reason inside \Box_T . Suppose we have $\Box_{T,x} A$. By cut-elimination we can find a T_x -proof p of A in which only formulas of complexity $< y$ occur. We now prove by induction on the subproofs of p , that all subconclusions of p are True_y . So A is True_y . Hence we find A. \square

6.2 A brief introduction to HA^* : In this section we describe the theory HA^* . This theory was introduced in Visser[82]. HA^* is to Beeson's fp-realizability (Beeson[75]) as Troelstra's $HA + \text{ECT}_0$ is to Kleene's r-realizability. This means that for a suitable variant of fp-realizability HA^* is the set of sentences such that their fp-translations are provable in HA. The natural way to define HA^* is by a fixed point construction as: HA plus the *Completeness Principle for HA^** . (Here it is essential that the construction is verifiable in HA, see below.) The Completeness Principle can be

viewed as an arithmetically interpreted modal principle. The Completeness Principle viewed modally is:

$$\mathfrak{C} \quad \vdash A \rightarrow \Box A$$

The Completeness Principle for a specific theory T is:

$$\mathfrak{C}[T] \quad \vdash A \rightarrow \Box_T A.$$

We have: $HA^* = HA + \mathfrak{C}[HA^*]$.

We briefly review some of the results of Visser[82] and DV[94].

- Let \mathbf{A} be the smallest class closed under atoms and all connectives except implication, satisfying: $A \in \Sigma_1$ and $B \in \mathbf{A} \Rightarrow (A \rightarrow B) \in \mathbf{A}$. Note that modulo provable equivalence in HA all formulas of the classical arithmetical hierarchy in their standard form are in \mathbf{A} . HA^* is conservative w.r.t. \mathbf{A} over HA. Note that $NNIL_0(\Sigma) \subseteq \mathbf{A}$.
- There are infinitely many incomparable T with $T = HA + \mathfrak{C}[T]$. However if $T = HA + \mathfrak{C}[T]$ *verifiably in HA*, then $T = HA^*$.
- Let KLS:=Kreisel-Lacombe-Shoenfield's Theorem on the continuity of the effective operations. We have: $HA^* \vdash KLS \rightarrow \Box_{HA^*} \perp$. This immediately gives Beeson's result that $HA \not\vdash KLS$ (see Beeson[75]).
- Every prime RE Heyting algebra \mathfrak{H} can be embedded into the Heyting algebra of HA^* . This mapping is primitive recursive in the enumeration of the generators of \mathfrak{H} w.r.t which \mathfrak{H} is RE. The mapping sends the generators to Σ -sentences. (See DV[94]).

We insert some basic facts on the provability logic of HA^* . These materials will be only needed in section 10. Clearly, HA^* satisfies the Löb conditions.

$$L1 \quad \vdash A \Rightarrow \vdash \Box A$$

$$L2 \quad \vdash \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$$

$$L3 \quad \vdash \Box A \rightarrow \Box \Box A$$

$$L4 \quad \vdash \Box(\Box A \rightarrow A) \rightarrow \Box A$$

i-K is given by IPC+L1,L2. i-L is i-K+L3,L4. We write i-K{P} for the extension of i-K with some principle P. We write ' $\mathfrak{C}[HA^*]$ ' in the modal contexts as C. Note that i-L{C} is valid for provability interpretations in HA^* .

A principle closely connected to C is the Strong Löb Principle:

$$SL \quad \vdash (\Box A \rightarrow A) \rightarrow A$$

As a special case of SL we have: $\vdash \neg \Box \perp$.

6.2.1 Fact: $i\text{-L}\{C\}$ is interderivable with $i\text{-K}\{SL\}$.

Proof: L4 is immediate from SL. “ $i\text{-K}\{SL\} \vdash C$ ”:

$$\begin{aligned} \vdash A &\rightarrow (\Box(A \wedge \Box A) \rightarrow (A \wedge \Box A)) \\ &\rightarrow A \wedge \Box A \\ &\rightarrow \Box A. \end{aligned}$$

“ $i\text{-L}\{C\} \vdash SL$ ”:

$$\begin{aligned} \vdash (\Box A \rightarrow A) &\rightarrow (\Box A \rightarrow A) \wedge \Box(\Box A \rightarrow A) \\ &\rightarrow (\Box A \rightarrow A) \wedge \Box A \\ &\rightarrow A. \end{aligned}$$

□

As a preliminary for section 10 we study the closed fragment of $i\text{-L}\{C\}$. A formula of the modal language is closed if it contains no propositional variables. We define: $\Box^0 \perp := \perp$, $\Box^{n+1} \perp := \Box \Box^n \perp$, $\Box^\omega \perp := \top$.

6.2.2 The Closed Fragment of $i\text{-L}\{C\}$: For every closed formula A there is an $\alpha =: \alpha^*(A)$, such that $\vdash A \leftrightarrow \Box^\alpha \perp$. It is easily seen that $\alpha^*(A)$ is unique.

Proof: The proof is by induction on A . We have:

$$\begin{aligned} \vdash \top &\leftrightarrow \Box^\omega \perp, \vdash \perp \leftrightarrow \Box^0 \perp \\ \vdash (\Box^\alpha \perp \wedge \Box^\beta \perp) &\leftrightarrow \Box^{\min(\alpha, \beta)} \perp \\ \vdash (\Box^\alpha \perp \vee \Box^\beta \perp) &\leftrightarrow \Box^{\max(\alpha, \beta)} \perp \\ \vdash (\Box^\alpha \perp \rightarrow \Box^\beta \perp) &\leftrightarrow \Box^{\alpha \rightarrow \beta} \perp, \text{ where } (\alpha \rightarrow \beta) := \top \text{ if } \alpha \leq \beta \text{ and } (\alpha \rightarrow \beta) := \beta \text{ if } \beta < \alpha. \end{aligned}$$

Note that $\min(\alpha, \gamma) \leq \beta \Leftrightarrow \gamma \leq (\alpha \rightarrow \beta)$.

$$\vdash \Box \Box^\alpha \perp \leftrightarrow \Box^{1+\alpha} \perp$$

□

7 Digression: assorted facts about Σ -preservativity: In this section we pause to look in some more detail at the notion of Σ -preservativity. We will transfer some classical results about Π -conservativity to Σ -preservativity and we formulate some principles of ‘preservativity logic’. *The only result of this section that is used in the rest of the paper is 7.1.*

Consider any consistent RE theory T , extending HA , in the language of HA . We introduce some notions, closely related to Σ -preservativity for T . Define (suppressing the index for T):

Provable Deductive Consequence	$A \triangleright_{\text{pdc}} B$	$:\leftrightarrow \Box(\Box A \rightarrow \Box B)$
Uniform Deductive Consequence	$A \triangleright_{\text{udc}} B$	$:\leftrightarrow \forall x \exists y (\Box_x A \rightarrow \Box_y B)$
Strong Uniform Deductive Consequence	$A \triangleright_{\text{sudc}} B$	$:\leftrightarrow \forall x (\Box_x A \rightarrow \Box_x B)$
Uniform Provably Deductive Consequence	$A \triangleright_{\text{updc}} B$	$:\leftrightarrow \forall x \exists y \Box(\Box_x A \rightarrow \Box_y B)$

We start with some equivalences and derivabilities.

7.1 Orey-Hájek for Σ -preservativity and Uniform Provable Deductive Consequence: T proves that the following are equivalent:

(i) $A \triangleright_{\Sigma} B$, (ii) $\forall x \Box(\Box_x A \rightarrow B)$, (iii) $A \triangleright_{\text{updc}} B$.

Proof: Reason in T . “(i) \rightarrow (ii)” Suppose $A \triangleright_{\Sigma} B$ and consider any x . We have $\Box(\Box_x A \rightarrow A)$ and $(\Box_x A) \in \Sigma$, hence $\Box(\Box_x A \rightarrow B)$.

“(ii) \rightarrow (iii)” From $\Box(\Box_x A \rightarrow B)$, we have that for some y $\Box_y(\Box_x A \rightarrow B)$ and hence $\Box\Box_y(\Box_x A \rightarrow B)$. Ergo: $\Box(\Box_y \Box_x A \rightarrow \Box_y B)$. Clearly for any u : $\Box(\Box_u A \rightarrow \Box_y \Box_u A)$. We may conclude: $\Box(\Box_x A \rightarrow \Box_y B)$. Hence: $A \triangleright_{\text{updc}} B$.

“(i) \leftarrow (iii)” Suppose $A \triangleright_{\text{updc}} B$ and $\Box(S \rightarrow A)$. It follows that for some x : $\Box\Box_x(S \rightarrow A)$. Moreover: $\Box(S \rightarrow \Box_x S)$. So $\Box(S \rightarrow \Box_x A)$. Hence for some y : $\Box(S \rightarrow \Box_y B)$. Ergo by reflection: $\Box(S \rightarrow B)$. \square

7.2 Fact: We have: $T \vdash A \triangleright_{\Sigma} B \rightarrow A \triangleright_{\text{udc}} B$.

Proof: Reason in T . Suppose $A \triangleright_{\Sigma} B$. Consider any x . For some y we have: $\Box_y(\Box_x A \rightarrow B)$. Suppose $\Box_x A$. Clearly for any u : $\Box_u A \rightarrow \Box_y \Box_u A$. It follows that $\Box_y B$. \square

7.3 Orey Hájek in HA^* : We have: $HA^* \vdash A \triangleright_{HA^*, \Sigma} B \leftrightarrow A \triangleright_{\text{udc}, HA^*} B$.

Proof: Immediate by the principle \mathfrak{C} , 7.1 and 7.2. \square

7.4 Fact: $HA^* \vdash (\Box_{HA^*} A \rightarrow B) \rightarrow A \triangleright_{HA^*, \Sigma} B$.

Proof:

$$\begin{aligned}
\vdash (\Box A \rightarrow B) &\rightarrow \forall x(\Box_x A \rightarrow B) \\
&\rightarrow \forall x \Box(\Box_x A \rightarrow B) \\
&\rightarrow A \triangleright_{\Sigma} B.
\end{aligned}$$

□

7.5 Modal logic for Σ -preservativity: Consider the language of modal propositional logic with a unary modal operator \Box and a binary operator \triangleright . We define arithmetical interpretations in the language of HA in the usual manner, interpreting \Box as provability in HA and \triangleright as Σ -preservativity over HA. We state the principles valid in HA for this logic, known at present. With the exception of $\Sigma 4$ and $\Sigma 8$ these principles are the duals of the principles of the interpretability logic ILM.

- 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 $\vdash \Box(\Box A \rightarrow A) \rightarrow \Box A$
 $\Sigma 1$ $\vdash \Box(A \rightarrow B) \rightarrow A \triangleright B$
 $\Sigma 2$ $\vdash A \triangleright B \wedge B \triangleright C \rightarrow A \triangleright C$
 $\Sigma 3$ $\vdash C \triangleright A \wedge C \triangleright B \rightarrow C \triangleright (A \wedge B)$
 $\Sigma 4$ $\vdash A \triangleright C \wedge B \triangleright C \rightarrow (A \vee B) \triangleright C$
 $\Sigma 5$ $\vdash A \triangleright B \rightarrow (\Box A \rightarrow \Box B)$
 $\Sigma 6$ $\vdash A \triangleright \Box A$
 $\Sigma 7$ $\vdash A \triangleright B \rightarrow (\Box C \rightarrow A) \triangleright (\Box C \rightarrow B)$
 $\Sigma 8$ Let X be a finite set of implications and let $Y := \{C \mid (C \rightarrow D) \in X\} \cup \{B\}$.
Take $A := \bigwedge X$. Then: $\vdash (A \rightarrow B) \triangleright \{A\}Y$

Here $(A \rightarrow B) \triangleright \{A\}Y$ is short for: $(A \rightarrow B) \triangleright \bigvee \{A\}Y$ and $\{C\}D$ is defined as follows:

- $\{C\}p := (C \rightarrow p)$, $\{C\}\perp := \perp$, $\{C\}\top := \top$,
- $\{C\}\Box E := \Box E$, $\{C\}(E \triangleright F) := (C \rightarrow (E \triangleright F))$, $\{C\}(E \rightarrow F) := (C \rightarrow (E \rightarrow F))$,
- $\{C\}(\cdot)$ commutes with \wedge and \vee .

The verification of L1-L4, $\Sigma 1$ - $\Sigma 3$, $\Sigma 4$ - $\Sigma 7$ is routine. For $\Sigma 4$ see 8.1 and for $\Sigma 8$ see 8.2.

From our principles e.g. Leivant's principle can be derived. (This is one of the Stellingen of Leivant[75].)

Le $\vdash \Box(A \vee B) \rightarrow \Box(A \vee \Box B)$

We leave this as an exercise to the reader.

It is open whether our axioms are complete.

7.6 Constructive versions of some results of Lindström and Švejdar

In his classical paper Švejdar[83], Vítěslav Švejdar studies logics for interpretability, conservativity and Rosser-orderings. Note that while *the study of methods* is Švejdar's, *the methods studied* are for a large part both Lindström's and Švejdar's.

The definitions of witness comparisons between sentences are as follows:

- $\exists x A x \leq_{\text{ex}} \exists y B y :\Leftrightarrow \exists x (A x \wedge \forall y < x \neg B y)$
- $\exists x A x <_{\text{ex}} \exists y B y :\Leftrightarrow \exists x (A x \wedge \forall y \leq x \neg B y).$

If B is decidable these definitions are constructively as useful as classically. However if we allow non-decidable B, the use of negation is rather heavy. I predict that more experimentation will reveal that the following notions are more useful and more pleasant to work with:

- $\exists x A x \leq_{\text{un}} \exists y B y :\Leftrightarrow \forall y (B y \rightarrow \exists x \leq y A x)$
- $\exists x A x <_{\text{un}} \exists y B y :\Leftrightarrow \forall y (B y \rightarrow \exists x < y A x).$

Classically we have:

$\exists x A x \leq_{\text{un}} \exists y B y$ is provably equivalent to the negation of $\exists x B x <_{\text{ex}} \exists y A y$,

$\exists x A x <_{\text{un}} \exists y B y$ is provably equivalent to the negation of $\exists x B x \leq_{\text{ex}} \exists y A y$,

$\exists x A x <_{\text{ex}} \exists y B y$ is provably equivalent to the negation of $\exists x B x \leq_{\text{un}} \exists y A y$,

$\exists x A x \leq_{\text{ex}} \exists y B y$ is provably equivalent to the negation of $\exists x B x <_{\text{un}} \exists y A y$.

Constructively all these connections fail.

We do have:

$$\text{O1} \quad \vdash \exists x A x \leq_{\text{ex}} \exists y B y \rightarrow \exists x A x \leq_{\text{un}} \exists y B y,$$

$$\text{O2} \quad \vdash \exists x A x <_{\text{ex}} \exists y B y \rightarrow \exists x A x <_{\text{un}} \exists y B y.$$

\leq_{ex} and $<_{\text{ex}}$ are provably transitive. We have by a 'downwards induction':

$$\text{O3} \quad \vdash \neg(\exists x A x <_{\text{ex}} \exists x A x).$$

But we do not even have: $\vdash \exists x A x \rightarrow \exists x A x \leq_{\text{ex}} \exists x A x$, since this expresses precisely the minimum principle for A, which is constructively not generally valid. (It is valid if A is decidable.)

We have:

$$\text{O4} \quad \leq_{\text{un}} \text{ provably satisfies the axioms of weak partial ordering,}$$

$$\text{O5} \quad <_{\text{un}} \text{ is provably transitive,}$$

$$\text{O6} \quad \vdash (\exists x A x <_{\text{un}} \exists x A x) \rightarrow \forall x \neg A x,$$

$$\text{O7} \quad \vdash \exists x A x <_{\text{un}} \exists y B y \rightarrow \exists x A x \leq_{\text{un}} \exists y B y,$$

$$\text{O8} \quad \vdash \exists x A x \leq_{\text{un}} \exists y B y <_{\text{un}} \exists z C z \rightarrow \exists x A x <_{\text{un}} \exists z C z,$$

$$\text{O9} \quad \vdash \exists x A x <_{\text{un}} \exists y B y \leq_{\text{un}} \exists z C z \rightarrow \exists x A x <_{\text{un}} \exists z C z.$$

Let $\Box^*A : \Leftrightarrow \exists x \Box_x A$. \Box_x is persistent, i.e. $\vdash (x < y \wedge \Box_x A) \rightarrow \Box_y A$. So we have:

$$\begin{aligned} \vdash \Box^*A <_{\text{ex}} \Box^*B &\Leftrightarrow \exists x (\Box_x A \wedge \neg \Box_x B), \\ \vdash \Box^*A \leq_{\text{un}} \Box^*B &\Leftrightarrow \forall x (\Box_x B \rightarrow \Box_x A) \Leftrightarrow B \triangleright_{\text{sudc}} A. \end{aligned}$$

7.6.1 Svejdar Principles: We give some constructive principles in the style of Švejdar [83]. The superiority of (iii) over (ii) is one of the arguments in favour of the universal witness comparison relation.

- i) $\vdash B \triangleright_{\Sigma} C \rightarrow A \triangleright_{\Sigma} ((\Box^*B \leq_{\text{un}} \Box^*A) \rightarrow C)$
- ii) $\vdash A \triangleright_{\Sigma} (\neg B \rightarrow \Box^*A <_{\text{ex}} \Box^*B)$
- iii) $\vdash A \triangleright_{\Sigma} ((B \rightarrow \Box^*A <_{\text{un}} \Box^*B) \rightarrow \Box^*A <_{\text{un}} \Box^*B)$
- iv) $\vdash A \triangleright_{\Sigma} ((B \rightarrow \Box^*A \leq_{\text{un}} \Box^*B) \rightarrow \Box^*A \leq_{\text{un}} \Box^*B)$

Proof: Reason in T

i) Suppose $B \triangleright_{\Sigma} C$. For any x :

$$\Box (\Box_x A \wedge \Box^*B \leq_{\text{un}} \Box^*A \rightarrow \Box_x B \rightarrow C).$$

$$\text{ii) } \Box (\Box_x A \wedge \neg B \rightarrow \Box_x A \wedge \neg \Box_x B \rightarrow \Box^*A <_{\text{ex}} \Box^*B).$$

iii) Reason in T. Consider x . Reason inside the \Box . Consider y . Suppose: $\Box_x A$, $(B \rightarrow \Box^*A <_{\text{un}} \Box^*B)$ and $\Box_y B$. We have to show: $\exists z < y \Box_z A$. If $y \leq x$, then $\Box_x B$ and hence B . So: $\Box^*A <_{\text{un}} \Box^*B$. We may conclude $\exists z < y \Box_z A$. If $y > x$, we get $\exists z < y \Box_z A$ from our assumption $\Box_x A$.

iv) By a minor modification of the reasoning under (iii). \square

7.6.2 Consequence: $\vdash A \triangleright_{\Sigma} (\Box^*B \leq_{\text{un}} \Box^*A \rightarrow B)$.

A version of the Feferman predicate can be defined as follows:

$$\bullet \Delta A : \Leftrightarrow \Box^*A <_{\text{ex}} \Box^* \perp.$$

Note that e.g. $\vdash \neg \Delta \perp$.

7.6.3 Fact: $\vdash A \triangleright_{\Sigma} \Delta A$.

Proof: Substitute \perp for B in 7.6.1(ii). \square

7.6.4 Orey Sentence: Let \triangleright^* be Π -conservativity or interpretability over T (in the classical context). A classical Orey sentence is a sentence G such that:

$$\top \triangleright^* G \text{ and } \top \triangleright^* \neg G.$$

The existence of an Orey-sentence shows the non-uniqueness of extension via interpretation or Π -conservative extension. Its existence also shows that the right hand side of Π -conservativity or interpretability for consistent T cannot be closed under conjunction. By duality an Orey sentence for Σ -preservativity (in the constructive context), should satisfy:

$$G \triangleright_{\Sigma} \perp \text{ and } \neg G \triangleright_{\Sigma} \perp.$$

Disanalogously it does *not* generally follow that the left hand side of Σ -preservativity is not closed under disjunction: in fact in the case of HA it is closed under disjunction. See 8.1.

Let G be such that $T \vdash G \leftrightarrow \neg \Delta G$. On the one hand: $G \triangleright_{\Sigma} \Delta G$ and hence $G \triangleright_{\Sigma} \neg G$. On the other hand $G \triangleright_{\Sigma} G$. Ergo $G \triangleright_{\Sigma} \perp$. Also: $\neg G \triangleright_{\Sigma} \Delta \neg G \triangleright_{\Sigma} \neg \Delta G \triangleright_{\Sigma} G$ and $\neg G \triangleright_{\Sigma} \neg G$. So $\neg G \triangleright_{\Sigma} \perp$.

In case $T=HA$, the results of section 8.1 give us: $(G \vee \neg G) \triangleright_{\Sigma} \perp$. \circ

7.6.5 Lindström's Theorem: Per Lindström proved that every equivalence class for interpretability or Π -conservativity contains both a Σ_2 - and a Π_2 -sentence. We explore his argument in our setting.

We first look at the Σ_2 -case. Pick by the Gödel Fixed Point Lemma (GFL) a Σ_2 -sentence R with $T \vdash R \leftrightarrow \Box^* B \leq_{ex} \Box^* R$, we have, using the fact that \triangleright_{Σ} is a semi-consequence relation for T :

- | | | |
|----|--|---------------|
| a) | $R \triangleright_{\Sigma} (\Box^* B \leq_{un} \Box^* R \rightarrow B)$ | 7.6.2 |
| b) | $R \triangleright_{\Sigma} (\Box^* B \leq_{ex} \Box^* R \rightarrow B)$ | a, O1, A1, A2 |
| c) | $R \triangleright_{\Sigma} \Box^* B \leq_{ex} \Box^* R$ | GFL, A1 |
| d) | $R \triangleright_{\Sigma} ((\Box^* B \leq_{ex} \Box^* R) \wedge (\Box^* B \leq_{ex} \Box^* R \rightarrow B))$ | b, c, A3 |
| e) | $R \triangleright_{\Sigma} B$ | d, A1, A2 |
| f) | $B \triangleright_{\Sigma} (\neg R \rightarrow \Box^* B \leq_{ex} \Box^* R)$ | 7.6.1(ii) |
| g) | $B \triangleright_{\Sigma} \neg \neg R$ | GFL, A1, A2 |

So we find: $R \triangleright_{\Sigma} B \triangleright_{\Sigma} \neg \neg R$. It is open whether we can find Q in Σ_2 , such that $B \equiv_{\Sigma} Q$.

We turn to the Π_2 -case. Pick by the Gödel Fixed Point Lemma a Π_2 -sentence R with $T \vdash R \leftrightarrow \Box^* B \leq_{un} \Box^* R$, we have:

- | | | |
|----|---|-------------|
| a) | $R \triangleright_{\Sigma} (\Box^* B \leq_{un} \Box^* R \rightarrow B)$ | 7.6.2 |
| b) | $R \triangleright_{\Sigma} \Box^* B \leq_{un} \Box^* R$ | GFL, A1 |
| c) | $R \triangleright_{\Sigma} ((\Box^* B \leq_{un} \Box^* R) \wedge (\Box^* B \leq_{un} \Box^* R \rightarrow B))$ | a, b, A3 |
| d) | $R \triangleright_{\Sigma} B$ | c, A1, A2 |
| e) | $B \triangleright_{\Sigma} ((R \rightarrow \Box^* B \leq_{un} \Box^* R) \rightarrow \Box^* B \leq_{un} \Box^* R)$ | 7.6.1(iv) |
| f) | $B \triangleright_{\Sigma} R$ | GFL, A1, A2 |
| g) | $B \equiv_{\Sigma} R$ | d, f |

So every equivalence class contains a Π_2 -sentence. \circ

7.6.6 Complexity of Σ -preservativity: Let U be the theory axiomatized by the set of A such that $T \vdash A \neg \neg$. Then U is a consistent theory extending PA in the language of PA. \triangleright^* stands for: conservativity. Let P_0 be a Π_1 -sentence, say P_0 is equivalent to $\neg S_0$, for $S_0 \in \Sigma_1$. We have:

$$\begin{aligned}
P_0 \triangleright_{T, \Sigma} \perp &\Leftrightarrow \forall S \in \Sigma_1 (\Box_T(S \rightarrow P_0) \Rightarrow \Box_T \neg S) \\
&\Leftrightarrow \forall P \in \Pi_1 (\Box_T(S_0 \rightarrow P) \Rightarrow \Box_T P) \\
&\Leftrightarrow \forall P \in \Pi_1 (\Box_U(S_0 \rightarrow P) \Rightarrow \Box_U P) \\
&\Leftrightarrow T \triangleright^*_{U, \Pi} S_0.
\end{aligned}$$

By a result, due independently to Per Lindström and Robert Solovay, the set $\{S_0 \in \Sigma_1 \mid T \triangleright^*_{U, \Pi} S_0\}$ is complete Π_2 . It follows that $\{P_0 \in \Pi_1 \mid P_0 \triangleright_{T, \Sigma} \perp\}$ is complete Π_2 . \circ

8 Closure properties of Σ -preservativity over HA: In this section we verify two closure properties of Σ -preservativity.

8.1 Closure under B1: We will show that Σ -preservativity is closed under B1.

We produce both proofs known to us. The first employs q-realizability. This form of realizability is a translation from \mathfrak{A} to \mathfrak{A} , due to Kleene. It is defined as follows:

- $xqP := P$, for P atomic
- $xq(A \wedge B) := ((x)_0 qA \wedge (x)_1 qB)$
- $xq(A \vee B) := (((x)_0 = 0 \rightarrow (x)_1 qA) \wedge ((x)_0 \neq 0 \rightarrow (x)_1 qB))$
- $xq(A \rightarrow B) := \forall y (yqA \rightarrow \exists z (\{x\}y \dot{=} z \wedge zqB)) \wedge (A \rightarrow B)$
- $xq\exists y A(y) := (x)_0 qA((x)_1)$
- $xq\forall y A(y) := \forall y \exists z (\{x\}y \dot{=} z \wedge zqA(y))$

The following facts can be verified in HA:

- a) $HA \vdash xqA \rightarrow A$,
- b) for every $A \in \Sigma$ and $\{y_1, \dots, y_n\}$ with $FV(A) \subseteq \{y_1, \dots, y_n\}$ we can find an index e such that: $HA \vdash A \rightarrow \exists z (\{e\}(y_1, \dots, y_n) \equiv z \wedge zqA)$.
- c) Suppose B_1, \dots, B_n, C have free variables among $\{y_1, \dots, y_m\}$ and that $\{x_1, \dots, x_n\}$ is disjoint from $\{y_1, \dots, y_m\}$. Then:
 $B_1, \dots, B_n \vdash_{HA} C \Rightarrow \exists e [x_1qB_1, \dots, x_nqB_n \vdash_{HA} \exists z (\{e\}(x_1, \dots, x_n, y_1, \dots, y_m) \equiv z \wedge zqC)]$.

The proofs are all simple inductions. For details the reader is referred to Troelstra[73], 188-202.

We will now show that Σ -preservativity satisfies B1. As is easily seen our proof is verifiable in HA.

We reason as follows:

$\alpha)$	$A \triangleright_{HA, \Sigma} C$	Assumption
$\beta)$	$B \triangleright_{HA, \Sigma} C$	Assumption
$\gamma)$	$S \in \Sigma$ and $HA \vdash S \rightarrow (A \vee B)$	Assumption.
$\delta)$	$xqS \vdash_{HA} \exists z (\{e\}(x) \equiv z \wedge zq(A \vee B))$	γ, e provided by (c)
$\epsilon)$	$S \vdash_{HA} \exists z (\{e\}(f) \equiv z \wedge zq(A \vee B))$	δ, f provided by (b)
$\eta)$	$S \wedge \exists z (\{e\}(f) \equiv z \wedge (z)_0 = 0) \vdash_{HA} A$	ϵ
$\zeta)$	$S \wedge \exists z (\{e\}(f) \equiv z \wedge (z)_0 \neq 0) \vdash_{HA} B$	ϵ
$\vartheta)$	$S \wedge \exists z (\{e\}(f) \equiv z \wedge (z)_0 = 0) \vdash_{HA} C$	η, α
$\iota)$	$S \wedge \exists z (\{e\}(f) \equiv z \wedge (z)_0 \neq 0) \vdash_{HA} C$	ζ, β
$\kappa)$	$S \vdash_{HA} \exists z (\{e\}(f) \equiv z \wedge (z)_0 = 0) \vee \exists z (\{e\}(f) \equiv z \wedge (z)_0 \neq 0)$	ϵ
$\lambda)$	$S \vdash_{HA} C$	$\vartheta, \iota, \kappa \quad \square$

The second proof employs a translation due to Dick de Jongh. For later use our definition is slightly more general than is really needed for the problem at hand. Let C be an \mathfrak{A} -sentence and let n be a natural number. Define a translation $[C]_n(\cdot)$ as follows:

- $[C]_nP := P$ for P atomic,
- $[C]_n(\cdot)$ commutes \wedge, \vee, \exists ,
- $[C]_n(A \rightarrow B) := ([C]_nA \rightarrow [C]_nB) \wedge \Box_n(C \rightarrow (A \rightarrow B))$,
- $[C]_n\forall y A(y) := \forall y [C]_nA(y) \wedge \Box_n(C \rightarrow \forall y A(y))$.

Let's first make a few quick observations, that make life easier:

- i) $HA \vdash [C]_nA \rightarrow \Box_n(C \rightarrow A)$
- ii) $HA \vdash [C]_n((A \rightarrow B) \wedge (A' \rightarrow B')) \leftrightarrow$
 $([C]_nA \rightarrow [C]_nB) \wedge ([C]_nA' \rightarrow [C]_nB') \wedge \Box_n(C \rightarrow ((A \rightarrow B) \wedge (A' \rightarrow B')))$.

Similarly for conjunctions of more than two implications.

$$\text{iii) } HA \vdash [C]_n \forall y \forall z A(y, z) \leftrightarrow (\forall y \forall z [C]_n A(y, z) \wedge \Box_n (C \rightarrow \forall y \forall z A(y, z))).$$

Similar for larger blocks of universal quantifiers.

$$\text{iv) } HA \vdash [C]_n \forall y (A(y) \rightarrow B(y)) \leftrightarrow \forall y ([C]_n A(y) \rightarrow [C]_n B(y)) \wedge \Box_n (C \rightarrow \forall y (A(y) \rightarrow B(y))).$$

$$\text{v) } HA \vdash [C]_n \forall y < z A(y) \leftrightarrow \forall y < z [C]_n A(y).$$

$$\text{vi) } \text{For } S \in \Sigma: HA \vdash S \leftrightarrow [C]_n S$$

(v) is immediate from the well know fact that: $HA \vdash \forall y < z \Box_n A(y) \rightarrow \Box_n \forall y < z A(y)$.

(vi) is immediate from (v).

Let's write $[C]_n \Gamma := \{[C]_n D \mid D \in \Gamma\}$. We have:

$$\text{vii) } \Gamma \vdash_{HA, n} A \Rightarrow [C]_n \Gamma \vdash_{HA} [C]_n A \text{ (verifiably in HA).}$$

Proof of (vii): The proof is by induction on the proof witnessing $\Gamma \vdash_{HA, n} A$. We treat two cases.

$\Gamma = \emptyset$ and A is an induction axiom, say for $B(x)$, of HA_n . Clearly $[C]_n A$ is HA-provably equivalent to:

$$[([C]_n B(0) \wedge \forall x ([C]_n B(x) \rightarrow [C]_n B(x+1)) \wedge \Box_n (C \rightarrow \forall x (B(x) \rightarrow B(x+1)))) \rightarrow (\forall x [C]_n B(x) \wedge \Box_n (C \rightarrow \forall x B(x)))] \wedge \Box_n (C \rightarrow A).$$

We have:

$$HA \vdash \Box_n A, \text{ and hence } HA \vdash \Box_n (C \rightarrow A).$$

So it follows that:

$$HA \vdash \Box_n (C \rightarrow \forall x (B(x) \rightarrow B(x+1))) \rightarrow \Box_n (C \rightarrow \forall x B(x)).$$

Moreover (as an instance of induction for $[C]_n B(x)$):

$$HA \vdash ([C]_n B(0) \wedge \forall x ([C]_n B(x) \rightarrow [C]_n B(x+1))) \rightarrow \forall x [C]_n B(x).$$

Combining these we find the promised: $HA \vdash [C]_n A$.

Suppose $A = (D \rightarrow E)$ and the last step in the proof was by:

$$\Gamma, D \vdash_{HA, n} E \Rightarrow \Gamma \vdash_{HA, n} D \rightarrow E.$$

From $\Gamma, D \vdash_{HA, n} E$, we have by the Induction Hypothesis: $[C]_n \Gamma, [C]_n D \vdash_{HA, n} [C]_n E$ and hence: $[C]_n \Gamma \vdash_{HA, n} [C]_n D \rightarrow [C]_n E$. Moreover for some finite $\Gamma_0 \subseteq \Gamma$, we have: $\Gamma_0, D \vdash_{HA, n} E$. Let B be the conjunction of the elements of Γ_0 . We find:

$$[C]_n \Gamma \vdash_{HA, n} \Box_n (C \rightarrow B) \text{ and } \vdash_{HA, n} \Box_n (B \rightarrow (D \rightarrow E)).$$

Hence: $[C]_n \Gamma \vdash_{HA, n} \Box_n (C \rightarrow (D \rightarrow E))$. We may conclude:

$$[C]_n \Gamma \vdash_{HA, n} ([C]_n D \rightarrow [C]_n E) \wedge \Box_n (C \rightarrow (D \rightarrow E)) \quad \square(\text{vii})$$

We now prove our principle. As is easily seen the argument can be verified in HA.

$\alpha)$	$A \triangleright_{HA, \Sigma} C$	Assumption
$\beta)$	$B \triangleright_{HA, \Sigma} C$	Assumption
$\gamma)$	$S \in \Sigma$ and $HA \vdash S \rightarrow (A \vee B)$	Assumption.
$\delta)$	for some n $S \vdash_{HA, n} (A \vee B)$	γ
$\epsilon)$	$[T]_n S \vdash_{HA} ([T]_n A \vee [T]_n B)$	$\delta, (vii)$
$\eta)$	$S \vdash_{HA} (\Box_n A \vee \Box_n B)$	$\epsilon, (vi), (i)$
$\zeta)$	$HA \vdash S \rightarrow C$	$\alpha, \beta, \eta, 7.1 \quad \square$

8.2 A closure rule for implication: To formulate our next closure rule it is convenient to work in a conservative extension of HA. Let \mathfrak{U}^+ be \mathfrak{U} extended with new predicate symbols (including the 0-ary case) for Σ -formulas. Let \mathfrak{f} be some assignment of Σ -formulas of \mathfrak{U} of the appropriate arities to the new predicate symbols.

Define: $[.]^\circ_n(.)$ and $\{.\}(.): \mathfrak{U}^+ \rightarrow \mathfrak{U}$ as follows:

- $\{A\}P := [A]^\circ_n P := P[\mathfrak{f}]$ for P atomic,
- $\{A\}(.)$ and $[A]^\circ_n(.)$ commute with \wedge, \vee and \exists ,
- $\{A\}(B \rightarrow C) := (A \rightarrow (B \rightarrow C))[\mathfrak{f}]$, $[A]^\circ_n(B \rightarrow C) := \Box_n((A \rightarrow (B \rightarrow C))[\mathfrak{f}])$,
- $\{A\}\forall x B(x) := (A \rightarrow \forall x B(x))[\mathfrak{f}]$, $[A]^\circ_n \forall x B(x) := \Box_n((A \rightarrow \forall x B(x))[\mathfrak{f}])$.

We have for B in \mathfrak{U}^+ (in the numbering of principles for $[.]$ of 8.1):

viii) $HA \vdash [A[\mathfrak{f}]]_n(B[\mathfrak{f}]) \rightarrow [A]^\circ_n B \rightarrow \{A\}B$,

ix) $[A]^\circ_n B$ is provably equivalent to a Σ -formula.

The proof of (viii) is an easy induction on B . (ix) is trivial.

8.2.1 Theorem: Suppose X is a set of implications in \mathfrak{U}^+ and let B be in \mathfrak{U}^+ . Say, $A := \bigwedge X$ and let $Y := \{C \mid (C \rightarrow D) \in X\} \cup \{B\}$. We have, verifiably in HA:

$$(A \rightarrow B)[\mathfrak{f}] \triangleright_{HA, \Sigma} \bigvee \{A\}Y.$$

Note that $\Sigma 8$ of section 7 is an immediate consequence of our principle.

Proof: To avoid heavy notations we suppress ' $[\mathfrak{f}]$ '. In the context of HA we assume that an \mathfrak{U}^+ -formula is automatically translated via \mathfrak{f} to the corresponding \mathfrak{U} -formula. Let S be a Σ -sentence (of \mathfrak{U}). Suppose $S \vdash_{HA} (A \rightarrow B)$. It follows that for some n $S \vdash_{HA, n} (A \rightarrow B)$ and hence by (vii) $[A]_n S \vdash_{HA} [A]_n (A \rightarrow B)$. By (ii) and (vi) we find:

$$S \vdash_{HA} (\bigwedge \{[A]_n C \rightarrow [A]_n D \mid (C \rightarrow D) \in X\} \wedge \Box_n (A \rightarrow A)) \rightarrow [A]_n B,$$

and so:

$$S \vdash_{HA} \bigwedge \{[A]_n C \rightarrow [A]_n D \mid (C \rightarrow D) \in X\} \rightarrow [A]_n B.$$

It follows by (viii) that:

$$S \vdash_{HA} \bigwedge \{ [A]_n^\circ C \rightarrow [A]_n D \mid (C \rightarrow D) \in X \} \rightarrow [A]_n^\circ B.$$

Since HA is a subtheory of PA, we get:

$$S \vdash_{PA} \bigwedge \{ [A]_n^\circ C \rightarrow [A]_n D \mid (C \rightarrow D) \in X \} \rightarrow [A]_n^\circ B.$$

By classical logic, we get: $S \vdash_{PA} \bigvee [A]_n^\circ Y$. Remember that PA is, verifiably in HA, conservative over HA w.r.t Π_2 -sentences (see Friedman[77]). Since $(S \rightarrow \bigvee [A]_n^\circ Y)$ is Π_2 , we get: $S \vdash_{HA} \bigvee [A]_n^\circ Y$. Ergo by (viii): $S \vdash_{HA} \bigvee \{A\} Y$. \square

Before closing this section we insert some remarks on the proof.

8.2.2 Remarks on the proof

- i) The above proof was obtained after analyzing an argument in De Jongh[82].
- ii) The step involving conservativity of PA over HA uses the Gödel-Friedman translation. Closer inspection reveals that the argument at hand just requires the Friedman translation. We give a sketch in 8.2.3. It follows that our results generalize to all essentially reflexive RE extensions T of HA that are closed both under the Friedman and the de Jongh translation.
- iii) From a sufficiently abstract perspective it would become clear that our present proof is just a variant of the proof of 4.2. $\{G \mid \Box_n(A \rightarrow G)\}$ in the present proof corresponds to the grey part of the Kripke model in the picture below 4.2. $[A]_n(\cdot)$ corresponds to the operation of adding the bottom-node in the picture (via Smoryński's operation) to the grey part. The detour via PA corresponds to the fact that our Kripke model argument is essentially classical. The special behaviour of Σ -sentences under translation corresponds with the fact that whether a Σ -sentence is forced or not (in a model of HA) is dependent just on the node under consideration and not on other nodes.

8.2.3 Appendix to 8.2: We show that the double negation translation can be eliminated from the proof of 8.2.1. We pick up the proof from the point where we have proved:

$$a) \quad S \vdash_{HA} \bigwedge \{ [A]_n^\circ C \rightarrow [A]_n D \mid (C \rightarrow D) \in X \} \rightarrow [A]_n^\circ B.$$

Let $H := \bigvee [A]_n^\circ Y$.

The Friedman translation $(E)^H$ of an arithmetical formula E is (modulo some details to avoid variable-clashes) the result of replacing all atomic formulas P in E by $(P \vee H)$. One easily shows:

- b) $HA \vdash E \Rightarrow HA \vdash (E)^H$;
- c) $HA \vdash H \rightarrow (E)^H$;
- d) for $S \in \Sigma$: $HA \vdash (S)^H \leftrightarrow (S \vee H)$.

By (a), (b) and (d) we have:

$$e) \quad S \vee H \vdash_{HA} \bigwedge \{ ((([A]^\circ_n C) \vee H) \rightarrow ([A]_n D)^H \mid (C \rightarrow D) \in X) \rightarrow ((([A]^\circ_n B) \vee H) \vee H) \}.$$

By (e) and propositional logic we find:

$$f) \quad S \vdash_{HA} \bigwedge \{ H \rightarrow ([A]_n D)^H \mid (C \rightarrow D) \in X \} \rightarrow H.$$

So by (f) and (c):

$$g) \quad S \vdash_{HA} H.$$

Hence by (g) and (viii): $S \vdash_{HA \vee \{A\} Y}.$ \square

9 On Σ -substitutions: In this section we prove our main results on Σ -substitutions. Given the results we already have, this is rather easy.

9.1 Fact: Let $f \in \Sigma^{\mathfrak{B}}$. We have:

- i) $\triangleright_{HA, \Sigma, f}$ is a σ -relation.
- ii) $\triangleright_{HA, \Sigma, \Sigma}$ is a σ -relation.

Proof: (ii) is a direct consequence of (i). We have:

- a) Every preservativity relation, and hence $\triangleright_{HA, \Sigma, f}$, satisfies A1-3.
- b) $\triangleright_{HA, \Sigma, f}$ satisfies B1 in virtue of 8.1
- c) $\triangleright_{HA, \Sigma, f}$ satisfies B2 by 3.5(iv)
- d) $\triangleright_{HA, \Sigma, f}$ satisfies B3 by 8.2, noting that $\{.\}(.)$ of 8.2 behaves like $.$ of section 3 on $B(\Sigma)$. \square

9.2 Theorem

- ii) $\triangleright_{ROB} = \triangleright_{HA, \Sigma, \Sigma} = \triangleright_{HA, \tau, \Sigma}.$

Proof: By 5.2(iii) \triangleright_{ROB} is the minimal σ -relation. So by 9.1(ii): $\triangleright_{ROB} \subseteq \triangleright_{HA, \Sigma, \Sigma}$. By 3.5(vi): $\triangleright_{HA, \Sigma, \Sigma} \subseteq \triangleright_{HA, \tau, \Sigma}$. We show $\triangleright_{HA, \tau, \Sigma} \subseteq \triangleright_{ROB}$. Suppose *not* $A \triangleright_{ROB} B$. Then by 5.2(i): $A^* \not\vdash_{IPC} B$. A^* is IPC-provably equivalent to a disjunction of prime NNIL-formulas. It follows that for some such disjunct, say C : $C \not\vdash_{IPC} B$. The Heyting algebra \mathfrak{H} axiomatized by C is prime and RE. By the embedding theorem proved in DV[94] (see section 6.2 of this paper) there is an $f \in \Sigma^{\mathfrak{B}}$, such that $HA^* \vdash C[f]$ and $HA^* \not\vdash B[f]$. Since $C[f] \in NNIL(\Sigma)$, we have by the $NNIL(\Sigma)$ -conservativity of HA^* over HA (proved in Visser[82]; see also section 6.2 for the statement of the full result), that $HA \vdash C[f]$. Since HA is a sub-theory of HA^* , we have: $HA \not\vdash B[f]$. Since $IPC \vdash C \rightarrow A$, it follows that $HA \vdash A[f]$ and $HA \not\vdash B[f]$. We may conclude: *not* $A \triangleright_{HA, \tau, \Sigma} B$. \square

10 The closed fragment of the provability logic of HA

Let IPC^{up} be the theory in language $\mathcal{L} := \mathcal{L}(\{\Box^n \perp \mid 0 \leq n \leq \omega\})$ axiomatized by IPC plus:

$$\Box^n \perp \vdash \Box^{n+1} \perp.$$

We identify \perp with $\Box^0 \perp$ and \top with $\Box^\omega \perp$. Let $\mathfrak{B} := \{\Box^n \perp \mid 0 \leq n \leq \omega\}$. The standard interpretation u of \mathcal{L} in \mathfrak{U} is the interpretation mapping $\Box^n \perp$ to $\Box_{\text{HA}}^n \perp$.

We will in this section notationally confuse \Box and \Box_{HA} . We write \Box^* for \Box_{HA}^* .

We reason as follows:

- | | | |
|---|--|--|
| a | for every $A \in \text{NNIL}(\Sigma)$ $\text{HA} \vdash \Box A \leftrightarrow \Box^* A$ | Visser[85] |
| b | for every $B \in B(\Sigma)$ $\text{HA} \vdash \Box B \leftrightarrow \Box B^*$ | section 9 |
| c | for every $C \in B(\Sigma)$ $\text{HA} \vdash \Box C \leftrightarrow \Box^* C^*$ | a, b ($C^* \in \text{NNIL}(\Sigma)$) |
| d | for every $\alpha \in \omega \cup \{\omega\}$ $\text{HA} \vdash \Box^\alpha \perp \leftrightarrow \Box^* \alpha \perp$ | a |
| e | for every $D \in B(\mathfrak{B})$ $\text{HA} \vdash \Box^*(D \leftrightarrow \Box^* \alpha^*(D) \perp)$ | d, section 6.2 |
| f | for every $D \in B(\mathfrak{B})$ $\text{HA} \vdash \Box^* D \leftrightarrow \Box^{1+\alpha^*(D)} \perp$ | e, a |
| g | for every $D \in B(\mathfrak{B})$ $\text{HA} \vdash \Box D \leftrightarrow \Box^{1+\alpha^*(D^*)} \perp$ | f, c |
| h | for every $E \in \mathcal{L}_{\Box}(\emptyset)$ $\exists \beta \text{HA} \vdash \Box E \leftrightarrow \Box^\beta \perp$ | g |

The last step is by induction on the box-depth of E . Note that we can read off from our proof an algorithm to compute β from E .

10.1 Sample computations

i) Consider $A := \neg \Box \perp \rightarrow \Box \Box \perp$. $A^* = \Box \perp \vee \Box \Box \perp$, which is equivalent to $\Box \Box \perp$. Hence:

$$\text{HA} \vdash \Box A \leftrightarrow \Box^3 \perp.$$

ii) Consider $A = (\neg \neg \Box \Box \perp \rightarrow \Box \Box \perp) \rightarrow (\Box \perp \vee \neg \Box \perp)$. Let $B := \neg \neg \Box \Box \perp \rightarrow \Box \Box \perp$. We compute A^* .

$$\begin{aligned} A &\equiv ((\Box \Box \perp \rightarrow (\Box \perp \vee \neg \Box \perp)) \wedge ([B](\neg \neg \Box \Box \perp) \vee ([B]\Box \perp) \vee ([B]\neg \Box \perp))) \\ &\equiv ((\Box \Box \perp \rightarrow (\Box \perp \vee \neg \Box \perp)) \wedge (\neg \neg \Box \Box \perp \vee \Box \perp \vee \neg \Box \perp)) \\ &\equiv ((\Box \Box \perp \rightarrow (\Box \perp \vee \neg \Box \perp)) \wedge (\Box \Box \perp \vee \Box \perp \vee \neg \Box \perp)). \end{aligned}$$

Clearly our last formula is equivalent to $\Box \perp \vee \neg \Box \perp$, so via the HA^* route we find:

$$\text{HA} \vdash \Box A \leftrightarrow \Box^2 \perp.$$

iii) Consider $A = (\neg \Box \perp \rightarrow \Box \perp) \rightarrow (\Box \Box \perp \vee \neg \Box \Box \perp)$. Let $B := \neg \Box \perp \rightarrow \Box \perp$. We compute A^* , using some shortcuts involving IPC^{up} -principles.

$$\begin{aligned} A &\equiv ((\Box \perp \rightarrow (\Box \Box \perp \vee \neg \Box \Box \perp)) \wedge ([B](\neg \Box \perp) \vee ([B]\Box \perp) \vee ([B]\neg \Box \perp))) \\ &\equiv (\neg \Box \perp \vee \Box \Box \perp \vee \neg \Box \Box \perp) \\ &\equiv (\Box \perp \vee \Box \Box \perp \vee \neg \Box \Box \perp). \end{aligned}$$

Clearly our last formula is equivalent to $\Box \Box \perp \vee \neg \Box \Box \perp$, so via the HA^* route we find:

$$\text{HA} \vdash \Box A \leftrightarrow \Box^3 \perp.$$

○

Define for $A \in \mathcal{U}$: $\alpha(A) := \alpha^*(A^*)$. We show that $\alpha(A)$ has a clear meaning in terms of the propositional theory IPC^{up} : it is the maximum of the α such that:

$$\text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A.$$

10.2 Fact: For $\alpha \in \{0, \dots, \omega\}$ and $A \in \mathcal{U}$:

$$\text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A \Leftrightarrow \alpha \leq \alpha(A).$$

Proof: We have:

$$\text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A \Rightarrow B, \Box^\alpha \perp \vdash_{\text{IPC}} A, \text{ for some conjunction } B \text{ of formulas of the form } \Box^\beta \perp \rightarrow \Box^{\beta+1} \perp.$$

Since B is in NNIL , we find:

$$B, \Box^\alpha \perp \vdash_{\text{IPC}} A \Leftrightarrow B, \Box^\alpha \perp \vdash_{\text{IPC}} A^*.$$

Hence:

$$\text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A \Leftrightarrow \text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A^*.$$

It follows that it is sufficient to prove:

$$\text{For } \alpha \in \{0, \dots, \omega\} \text{ and } A \in \text{NNIL}(\{\Box^n \perp \mid 1 \leq n\}):$$

$$\text{IPC}^{\text{up}} \vdash \Box^\alpha \perp \rightarrow A \Leftrightarrow \alpha \leq \alpha^*(A).$$

The proof is by induction on A . We leave the simple argument to the reader. The case of disjunction uses the fact that $\text{IPC}^{\text{up}} + \Box^\alpha \perp$ has the disjunction property. \square

References

- Beeson, M.J., 1975, *The nonderivability in intuitionistic formal systems of theorems on the continuity of effective operations*, JSL 40, 321-346.
- Berarducci, A., 1990, *The interpretability logic of Peano arithmetic*, JSL 55, 1059-1089.
- Boolos, G., 1993, *The logic of provability*, Cambridge University Press, Cambridge, UK.
- De Jongh, D.H.J., 1982, *Formulas of one propositional variable in intuitionistic arithmetic*, in: van TV[82], 51-64.
- De Jongh, D.H.J., Renardel, G.R., Van Benthem, J.F.A.K., Visser, A., 1994, *NNIL, a study in intuitionistic propositional logic*, Logic Group Preprint Series 111, Dept. of Philosophy, University of Utrecht, Heidelberglaan 8, 3584 CS Utrecht.
- De Jongh, D.H.J., Visser, A., 1994, *Embeddings of Heyting Algebras, revised version*, Logic Group Preprint Series 115, Dept. of Philosophy, University of Utrecht, Heidelberglaan 8, 3584 CS Utrecht.
- De Jongh, D.H.J., Chagrova, L.A., to appear, *The decidability of dependency in intuitionistic propositional logic*.
- Esakia, L., Artemov, S., 1991, *Provability logic*, Special issue of Studia Logica, L1.

- Fagin, R., Halpern, J.Y., Vardi, M.Y., 1992, *What is an inference rule?* JSL, vol 57, 1018-1045.
- Friedman, H., 1975, *One hundred and two problems in mathematical logic*, JSL 40, 113-129.
- Friedman, H., 1977, *Classically and intuitionistically provably recursive functions*, in: MS[77], 21-27.
- Leivant, D., 1975, *Absoluteness in intuitionistic logic*, PhD Thesis, University of Amsterdam. (Corrected reprint: 1979, Mathematical Centre Tract, no 73, Amsterdam.)
- Leivant, D., 1980, *Innocuous substitutions*, JSL 45, 363-368.
- Leivant, D., 1981, *Implicational complexity in intuitionistic arithmetic*, JSL 46, 240-248.
- Müller, G.H., and Scott, D. (eds.), 1977, *Higher set theory*, Springer Lecture Notes in Mathematics 669, Springer, Berlin.
- Petkov, P.P. (ed.), 1990, *Mathematical Logic*, (Proceedings of the Heyting Conference, Chaika, 1988), Plenum press, New York and London.
- Pitts, A., 1992, *On an interpretation of second order quantification in first order intuitionistic propositional logic*, JSL 57, 33-52.
- Renardel de Lavalette, G.R., 1986, *Interpolation in a fragment of intuitionistic propositional logic*, Logic Group Preprint Series 5, Dept. of Philosophy, University of Utrecht, Heidelberglaan 8, 3584 CS Utrecht.
- Rybakov, V.V., 1992, *Rules of inference with parameters for intuitionistic logic*, JSL 57, 912-923.
- Shavrukov, V., 1993, *Subalgebras of diagonalizable algebras of theories containing arithmetic*, Dissertationes Mathematicae, Polska Akademia Nauk., Mathematical Institute.
- Statman, R., 1979, *Intuitionistic propositional logic is polynomial space complete*, Theoretical Computer Science, vol. 9, 67-72.
- Švejdar, V., 1983, *Modal Analysis of Generalized Rosser Sentences*, JSL 48, 986-999.
- Smoryński, C., 1973, *Applications of Kripke Models*, in Troelstra[73], 324-391.
- Smoryński, C., 1985, *Self-Reference and Modal Logic*, Springer Verlag, Berlin
- Troelstra, A.S. (ed.), 1973, *Metamathematical Investigations of Intuitionistic Arithmetic and Analysis*, Springer Lecture Notes 344, Springer Verlag, Berlin.
- Troelstra, A.S. and Van Dalen, D (eds.), 1982, *The L.E.J. Brouwer Centenary Symposium*, North Holland, Amsterdam.
- Troelstra, A.S. and Van Dalen, D, 1988a, *Constructivism in Mathematics*, vol 1, North Holland, Amsterdam.

- Troelstra, A.S. and Van Dalen, D., 1988b, *Constructivism in Mathematics*, vol 2, North Holland, Amsterdam.
- Van Benthem, J.F.A.K., 1991, *Temporal Logic*, Report X-91-05, ILLC, University of Amsterdam; to appear in *Handbook of Logic in AI and Logic Programming* (Dov Gabbay et al., eds., Oxford University Press).
- Van Oosten, J., 1991, *Exercises in realizability*, PhD Thesis, University of Amsterdam.
- Visser, A., 1982, *On the Completeness Principle*, *Annals of Mathematical Logic* 22, 263-295.
- Visser, A., 1985, *Evaluation, provably deductive equivalence in Heyting's Arithmetic of substitution instances of propositional formulas*, Logic Group Preprint Series 4, Dept. of Philosophy, University of Utrecht, Heidelberglaan 8, 3584 CS Utrecht.
- Visser, A., 1990, *Interpretability Logic*, in: Petkov[90], 175-209.
- Zambella, D., 1994, *Shavrukov's Theorem on the subalgebras of diagonalizable algebras for theories containing IA_0+Exp* , *Notre Dame Journal of Formal Logic*, vol.35, 147-157.