

# A new method for establishing conservativity of classical systems over their intuitionistic version

THIERRY COQUAND<sup>†</sup> and MARTIN HOFMANN<sup>‡</sup>

<sup>†</sup> Chalmers University, Computer Science, S41296 Göteborg, Sweden.  
E-mail: coquand@cs.chalmers.se

<sup>‡</sup> Laboratory for Foundations of Computer Science, The King's Buildings, Mayfield Road, Edinburgh EH9 3JZ.  
E-mail: mxh@dcs.ed.ac.uk

Received 19 May 1998; revised 5 March 1999

We use a syntactical notion of Kripke models to obtain interpretations of subsystems of arithmetic in their intuitionistic counterparts. This yields, in particular, a new proof of Buss' result that the Skolem functions of Bounded Arithmetic are polynomial time computable.

## 1. Introduction

One goal of Hilbert's programme (Hilbert 1971) was to make sense of classical reasoning, which makes use of the law of excluded-middle, from an intuitionistic viewpoint. If so understood, it has been answered positively for the case of Peano arithmetic  $PA$  by the usual Gödel–Gentzen–Bernays negative translation of this system into its intuitionistic version  $HA$  (Heyting arithmetic). This translation is a powerful and purely syntactical method for reducing a classical system to its intuitionistic version. We remark that this result motivated the distinction between *intuitionism* and *finitism*, which is usually equated with primitive recursive arithmetic  $PRA$  (Troelstra and van Dalen 1988). It follows from the work of Gödel that there is no interpretation of  $PA$  or  $HA$  in  $PRA$ ; however Parsons' result, discussed below, shows that there are *a priori* non-finitary systems that can be interpreted in  $PRA$ .

Unfortunately, the negative translation does not always work.

The most prominent example is the system  $PA^\omega + AC$ , where  $PA^\omega$  is an extension of  $PA$  in which one can quantify over functions, functionals, ... and  $AC$  is the Axiom of Choice

$$(\forall x)(\exists y)A(x, y) \rightarrow (\exists f)(\forall x)A(x, f(x)).$$

In this case, the negative interpretation of  $AC$ , which has the form

$$(\forall x)\neg\neg(\exists y)A^*(x, y) \rightarrow \neg\neg(\exists f)(\forall x)A^*(x, f(x)),$$

cannot be proved in the intuitionistic version  $HA^\omega + AC$ . It can be shown, however, in this case that the proof-theoretic strength of the classical version  $PA^\omega + AC$  is strictly greater than the one of its intuitionistic version  $HA^\omega + AC$ . This excludes *a priori* any purely syntactical translation of the classical system into its intuitionistic version.

There are other cases, where it is known by other methods (such as ordinal analysis) that the two versions have the same strength, but where the negative translation nevertheless does not seem to work. An example of this latter case is the system  $I\Sigma_1^0 + EM$ , which is a subsystem of  $PA$  in which induction is restricted to  $\Sigma_1^0$  (existential) statements.

The negative translation of an existential statement  $(\exists x)\varphi$  for atomic  $\varphi$  is easily seen to be equivalent to  $\neg\neg(\exists x)\varphi$ , which is not, in general, an existential statement. But it seems that we need induction over such statements in order to interpret the negative translation of  $I\Sigma_1^0$ !

In this case, it is natural to ask if there is not a purely syntactical translation of the classical system in its intuitionistic version, which would avoid the use of an ordinal analysis.

In this paper we present such a method.

Applied to  $I\Sigma_1^0 + EM$ , it shows directly the conservativity of this system over its intuitionistic version for  $\Pi_2^0$  statements. Since it is a direct consequence of modified realisability that  $I\Sigma_1^0$  is conservative over primitive recursive arithmetic for  $\Pi_2^0$  statements, this can be seen as yet another proof of Parsons' result (Parsons 1970):  $I\Sigma_1^0 + EM$  is conservative over primitive recursive arithmetic for  $\Pi_2^0$  statements. This result is important since primitive recursive arithmetic is often thought to represent exactly finitism (as a restricted form of intuitionism). This conservativity establishes that one can make sense in a finitary way of classical logic with  $\Sigma_1^0$  induction, schema, which requires *a priori* some appeal to infinity (a set defined by a  $\Sigma_1^0$  formula is in general not decidable).

The method we present can be extended to an interpretation of systems with König's lemma and the axiom of choice. We only sketch this extension and instead apply in detail the method to  $CPV^\omega$  – a higher-order extension of Buss' Bounded Arithmetic introduced by Cook and Urquhart (Cook and Urquhart 1993). We obtain a new and relatively simple proof that the first-order functions definable therein are of polynomial complexity.

Our method is reminiscent of an argument due to Buchholz (Buchholz 1977), which he used to show that a positive inductive definition ( $ID_1$ ) can be translated into a strictly positive one.

## 2. Peano arithmetic with $\Sigma_1$ -induction

In this section we present a method to interpret a system equivalent to  $I\Sigma_1^0 + EM$ , that is, a fragment of Peano arithmetic with induction restricted to  $\Sigma_1^0$ -formulas, to its intuitionistic version. We emphasise that the result in itself is not new; a proof using cut elimination can be found in Parsons (1970). A proof using Gödel's *Dialectica* interpretation can be found in Parsons (1972) and Jervell (1998). The volume Aczel *et al.* (1992) contains two further different proofs (Wainer and Wallen 1992; Pohlers 1992). Indeed, it was the analogy between the proof in Wainer and Wallen (1992) and the  $\Omega$ -rule of Buchholz (Buchholz 1981; Buchholz 1977) that suggested the present method.

The advantage of proofs using functional interpretations, such as the present one or that involving the *Dialectica* interpretation, is that the extracted functional witness for a  $\Pi_2^0$  is obtained directly by a structural induction on proofs and does not involve an exhaustive search component as witnesses obtained via cut elimination do. One might

thus hope that our method (as well as the one based on the *Dialectica* interpretation) might have applications to program extraction from classical proofs. Our method is rather different from the *Dialectica* interpretation even at the level of the extracted programs. For example, the *Dialectica* interpretation translates an instance of the contraction rule as a case distinction, whereas under our interpretation it is interpreted as a duplication of a variable as in realisability.

Of course, whether or not this results in a gain of efficiency in the extracted programs remains to be found out by carrying out practical examples.

### 2.1. The system $PRA^\omega$

We consider a variation of the system  $HA^\omega$  (Troelstra and van Dalen 1988) called  $PRA^\omega$  where:

- the induction scheme is restricted to formulae of the form  $A(x) \stackrel{\text{def}}{=} \exists y.g(x, y)=0$  (these are called purely existential or  $\Sigma_1^0$ -formulae)
- the target type of primitive recursive definitions is restricted to ground type  $N$ .

We also include the axiom  $x=0 \vee \neg x=0$  expressing decidability of atomic formulas. It can be seen that in the presence of this axiom, quantifier-free formulas are equivalent to atomic formulas. On the other hand, with induction over  $\Sigma_1^0$  formulas with quantifier-free kernel, this would be provable.

It would be possible to use a restricted version of this system with variables ranging only over natural numbers and functions, like the system of elementary analysis  $EL_1$  used in Troelstra (1974).

By modified realisability, it can be seen that if  $PRA^\omega \vdash (\forall x)(\exists y) A(x, y)$ , there exists  $f : N \rightarrow N$  primitive recursive such that  $A(n, f(n))$  for all  $n$ .

The modified realisability we use is completely standard; see, for instance, Troelstra and van Dalen (1988). It will be spelt out below for the similar system  $IPV^\omega$ .

Our task is to interpret  $PRA^\omega + EM$  in  $PRA^\omega$ . As mentioned above, the negative translation in itself does not work because the translation does not preserve  $\Sigma_1^0$ -formulas. However, it does provide a syntactical translation of  $PRA^\omega + EM$  into the system  $PRA^\omega + MP$  where  $MP$  is *Markov's principle* (Troelstra and van Dalen 1988), which is the schema

$$\neg\neg(\exists x)\varphi \rightarrow (\exists x)\varphi$$

for each purely existential formula  $(\exists x)\varphi$ . What we present now, is an interpretation of  $PRA^\omega + MP$  into  $PRA^\omega$ . The idea is to use a Kripke model, internally definable in  $PRA^\omega$ , which will be a model of  $PRA^\omega + MP$ .

### 2.2. Kripke semantics for Markov's principle

Let  $s, t : N \rightarrow N$  be (not necessarily closed) terms. We write  $st$  for the term  $\lambda x:N.s(x)t(x)$  (point-wise multiplication) and  $T(s)$  for the formula  $\exists x.s(x)=0$ . The worlds of our Kripke model are terms  $t : N \rightarrow N$  thought of as codes for  $\Sigma_1^0$  formulae via  $T(-)$ . Accordingly, the partial order  $s \leq t$  is defined to be  $T(t) \rightarrow T(s)$ . Notice that for this partial order, there is a greatest element  $\lambda x:N.1$  and  $st$  is a greatest lower bound of  $s$  and  $t$ , that is,

$T(st) \leftrightarrow T(s) \vee T(t)$ . Thus, we have a meet-semi lattice. (In the case of  $ID_1$ , the meet-semi lattice can be formed by taking sequents that contain the inductively defined predicate only in positive positions; this is implicit in the reduction of  $ID_v^i + EM$  to  $ID_v^i$ , for the simplest case  $v = 1$ , in Buchholz (1981, pp. 224–227).)

To each formula  $A$  of  $PRA^\omega$ , we now associate another formula  $f \Vdash A$  with one extra variable  $f : N \rightarrow N$ . The defining clauses, which follow the usual definition of Kripke semantics, are as follows

$$\begin{aligned} f \Vdash A &\stackrel{\text{def}}{=} A \vee T(f) \quad \text{when } A \text{ is atomic} \\ f \Vdash A \wedge B &\stackrel{\text{def}}{=} f \Vdash A \wedge f \Vdash B \\ f \Vdash A \vee B &\stackrel{\text{def}}{=} f \Vdash A \vee f \Vdash B \\ f \Vdash \neg A &\stackrel{\text{def}}{=} (\forall g)[g \Vdash A \rightarrow T(fg)] \\ f \Vdash A \rightarrow B &\stackrel{\text{def}}{=} (\forall g)[g \Vdash A \rightarrow f g \Vdash B] \\ f \Vdash (\forall x)A &\stackrel{\text{def}}{=} (\forall x)f \Vdash A \\ f \Vdash (\exists x)A &\stackrel{\text{def}}{=} (\exists x)f \Vdash A. \end{aligned}$$

Notice that we could write the following alternatives for the clauses of implication and negation

$$\begin{aligned} f \Vdash \neg A &\stackrel{\text{def}}{=} (\forall g \leq f)[g \Vdash A \rightarrow T(g)] \\ f \Vdash A \rightarrow B &\stackrel{\text{def}}{=} (\forall g \leq f)[g \Vdash A \rightarrow g \Vdash B] \end{aligned}$$

and get an equivalent definition.

We remark that in an earlier version of this work, we used a composition of the present translation with the negative translation; we then get a direct interpretation of the law of excluded middle and the truth-values get a structure similar to the one described in Sambin (1996).

**Lemma 2.1.** If  $A$  is provable in  $PRA^\omega$  without using induction then  $f \Vdash A$  is provable in  $PRA^\omega$ .

*Proof.* The proof is immediate by induction on derivations in  $PRA^\omega$ . □

**Lemma 2.2.** If  $A$  is  $\Sigma_1^0$ , then  $f \Vdash A$  is equivalent (in  $PRA^\omega$ ) to  $A \vee T(f)$ .

*Proof.* In this case  $A$  is of the form  $(\exists n)[t(n) = 0]$ , so  $f \Vdash A$  is  $(\exists n)[t(n)=0 \vee T(f)]$ , which is equivalent to  $A \vee T(f)$ . □

**Lemma 2.3.** If  $A$  is  $\Sigma_1^0$ , then  $f \Vdash \neg A$  is equivalent (in  $PRA^\omega$ ) to  $A \rightarrow T(f)$ .

*Proof.*  $f \Vdash \neg A$  is  $(\forall g)[g \Vdash A \rightarrow T(fg)]$ . By Lemma 2.2,  $g \Vdash A$  is equivalent to  $A \vee T(g)$ , and hence  $h \Vdash \neg A$  is equivalent to  $A \rightarrow T(f)$ . □

**Lemma 2.4.** If  $A$  is  $\Sigma_1^0$ , then  $f \Vdash \neg\neg A$  is equivalent (in  $PRA^\omega$ ) to  $A \vee T(f)$ .

*Proof.* It is clear that  $f \Vdash A$ , which is equivalent to  $A \vee T(f)$  by Lemma 2.2, implies  $f \Vdash \neg\neg A$ . In the other direction,  $f \Vdash \neg\neg A$  is  $(\forall g)[g \Vdash \neg A \rightarrow T(fg)]$  while  $g \Vdash \neg A$  is equivalent to  $A \rightarrow T(g)$  by Lemma 2.3. If we take  $g$  such that  $A$  is equivalent to  $T(g)$ , we get that  $f \Vdash \neg\neg A$  implies  $T(fg)$ , which is equivalent to  $A \vee T(f)$ , as desired. □

**Corollary 2.5.** (Markov’s Principle) If  $A$  is  $\Sigma_1^0$ , then  $f \Vdash \neg\neg A \rightarrow A$ .

**Lemma 2.6.** If  $A(x)$  is  $\Sigma_1^0$ , then  $f \Vdash A(0) \rightarrow (\forall x)[A(x) \rightarrow A(Sx)] \rightarrow A(x)$  is provable in  $PRA^\omega$ .

*Proof.* The proof is direct by Lemma 2.2. □

All this gives a proof of

**Theorem 2.7.** If  $A$  is provable in  $PRA^\omega + MP$ , then  $f \Vdash A$  is provable in  $PRA^\omega$  for any  $f$ .

There are two possible ways of seeing that a  $\Sigma_1^0$  formula  $A$  provable in  $PRA^\omega + MP$  is provable in  $PRA^\omega$  using this theorem: we take  $f = \lambda n.1$  or  $f$  such that  $T(f)$  is equivalent to  $A$ . In both cases,  $f \Vdash A$  is equivalent to  $A$  (in  $PRA^\omega$ ), and hence by the theorem,  $A$  is provable in  $PRA^\omega$ . Using a negative translation, it follows that if a  $\Pi_2^0$ -formula  $A$  is provable in  $PRA^\omega + EM$ , then it is provable in  $PRA^\omega$ .

We can see that our method is quite similar to Friedman’s translation (Troelstra and van Dalen 1988); the important difference being that the disjunctively added formula is now a parameter that is ‘variable’ and gets changed during the translation.

### 3. Bounded arithmetic

In his thesis (Buss 1986) Sam Buss introduced a fragment of Peano-Arithmetic  $S_2^1$  in which the induction scheme is replaced by the weaker scheme of so-called *NP-induction*.

If  $A(x)$  is a formula with free variable  $x$ , let us define

$$PIND(A(x)) \stackrel{\text{def}}{=} A(0) \rightarrow (\forall x.A(\lfloor x/2 \rfloor) \rightarrow A(x)) \rightarrow \forall x.A(x).$$

The formula  $PIND(A(x))$  is called an instance of *polynomial induction*. An instance of *NP-induction* is a formula  $PIND(A(x))$  where  $A(x)$  is a bounded  $\Sigma_1^0$ -formula, i.e., one of the form

$$A(x) \stackrel{\text{def}}{=} \exists y \leq t(x).s(x, y)=0$$

where  $s, t$  are terms and  $t$  does not contain  $y$ .

In order that enough bounded  $\Sigma_1^0$ -formulas exist, one needs a fair number of basic functions and quantifier free axioms. We will discuss this point below for a certain extension of bounded arithmetic.

We also remark that in Buss’ formulation of  $S_2^1$  the kernel  $s(x, y) = 0$  is replaced by a more general concept (sharply bounded formula), which turns out to be equivalent in the more general system that we consider below.

The main result of Buss (1986) is that the Skolem functions of  $\Pi_0^2$ -statements that are provable in  $S_2^1$  are polynomial time computable (*PTIME*). Buss proves this result by assigning *PTIME*-functions to cut-free proofs. Cook and Urquhart (Cook and Urquhart 1993) give an alternative proof involving a *Dialectica* interpretation of  $S_2^1$  in  $IPV^\omega$  – a higher-order generalisation of *intuitionistic* bounded arithmetic. The intuitionistic system  $IPV^\omega$  admits a straightforward realisability interpretation in  $PV^\omega$  – a *PTIME*-variant of Gödel’s  $T$  from which the desired result follows directly.

We will now show how our method gives an alternative proof of Buss’ result. The discussion of advantages of functional interpretations compared to cut elimination from Section 2 applies to this case. Moreover, since our method is relatively simple, it might be possible to apply it to the weak subsystems of bounded arithmetic studied by Johannsen (Johannsen 1996). He reports that he did not succeed in applying the Dialectica translation whereas Buss’ method involving cut elimination does go through.

### 3.1. The systems $PV^\omega$ and $IPV^\omega$

The system  $PV^\omega$  is the simply-typed lambda calculus over one base type  $o$  (for natural numbers in binary notation) and constants with types as indicated.

- The constant zero:  $0 : o$ .
- The two successor functions:  $s_0, s_1 : o \rightarrow o$ .
- Integer division by two (‘mix-fix notation’):  $\lfloor \frac{\_}{2} \rfloor$  or  $\lfloor \_ / 2 \rfloor$ .
- The parity function *Parity* :  $o \rightarrow o$ .
- The (infix) functions *chop*, *pad*, and *smash* :  $\dashv, \boxplus, \# : o \rightarrow o \rightarrow o$
- The ternary conditional *Cond* :  $o \rightarrow o \rightarrow o \rightarrow o$ .
- The bounded recursor  $\mathcal{R} : o \rightarrow (o \rightarrow o \rightarrow o) \rightarrow (o \rightarrow o) \rightarrow o \rightarrow o$

The intended interpretation of the constants is as follows:  $0 = 0, s_0(x) = 2x, s_1(x) = 2x + 1, Parity(2x) = 0, Parity(2x + 1) = 1, x \dashv y = \lfloor x/2^{|y|} \rfloor, x \boxplus y = x \cdot 2^{|y|}, x \# y = 2^{|x| \cdot |y|}$  where  $|x| = \lceil \log_2(x + 1) \rceil$  is the length of  $x$  in binary notation. The meaning of the conditional is  $Cond(0, y, z) = y, Cond(x + 1, y, z) = z$ . The meaning of the recursor  $\mathcal{R}$  is given by

$$\mathcal{R}(g, h, k, x) = Cond(x, g, Cond(t \dashv k(x), t, k(x)))$$

where  $t \stackrel{\text{def}}{=} h(x, \mathcal{R}(g, h, k, \lfloor \frac{x}{2} \rfloor))$ . To understand this definition assume functions  $g, h, k$  of appropriate type and let  $f$  be the function defined by

$$\begin{aligned} f(0) &= g \\ f(x) &= h(x, f(\lfloor \frac{x}{2} \rfloor)) \text{ when } x > 0. \end{aligned}$$

We have  $\mathcal{R}(g, h, k, x) \leq 2^{|k(x)|}$  when  $x > 0$  and  $\mathcal{R}(g, h, k, x) = f(x)$  provided that  $f(x) \leq 2^{|k(x)|}$ , and thus, in particular, if  $f(x) \leq k(x)$ . Therefore, the recursor  $\mathcal{R}$  admits the definition of functions by Cobham’s scheme of *bounded recursion on notation* and it follows from Cobham’s theorem (Cobham 1965) that all *P*TIME -functions can be defined in  $PV^\omega$ .

The system  $IPV^\omega$  is a many-sorted intuitionistic predicate calculus over  $PV^\omega$ . Its non-logical axioms are the defining equations for the  $PV^\omega$ -constants (for a precise definition see Cook and Urquhart (1993)) and all formulas  $PIND(\exists y \leq t(x).s(x, y)=0)$  where  $t$  contains free variables of base type only. We also include the axiom  $x=0 \vee \neg x=0$  expressing decidability of atomic formulas. As in the case of  $PRA^\omega$ , it follows that quantifier-free formulas are equivalent to atomic formulas.

The bounded quantifier  $\exists x \leq t.A$  is shorthand for  $\exists x.Lesseq(x, t)=0 \wedge A$ , where *Lesseq* is a  $PV^\omega$  term denoting comparison of integers.

3.2. Realisability for  $IPV^\omega$

The definition of realisability for  $IPV^\omega$  is completely standard. To each  $IPV^\omega$ -formula  $A$ , we associate a formula  $\tilde{x} \textcircled{R} A$  by the following definition:

$$\begin{aligned}
 () \textcircled{R} A & \stackrel{\text{def}}{=} A, \text{ if } A \text{ is atomic} \\
 \tilde{x}, \tilde{y} \textcircled{R} A \wedge B & \stackrel{\text{def}}{=} \tilde{x} \textcircled{R} A \wedge \tilde{y} \textcircled{R} B \\
 z, \tilde{x}, \tilde{y} \textcircled{R} A \vee B & \stackrel{\text{def}}{=} (z=0 \wedge \tilde{x} \textcircled{R} A) \vee (z=1 \wedge \tilde{x} \textcircled{R} B) \\
 \tilde{y} \textcircled{R} A \rightarrow B & \stackrel{\text{def}}{=} (\forall \tilde{x}. \tilde{x} \textcircled{R} A \rightarrow \tilde{y}(\tilde{x}) \textcircled{R} B) \\
 \tilde{x} \textcircled{R} \forall y. A & \stackrel{\text{def}}{=} \forall y. \tilde{x}(y) \textcircled{R} A \\
 z, \tilde{x} \textcircled{R} \exists y. A & \stackrel{\text{def}}{=} \tilde{x} \textcircled{R} A(z).
 \end{aligned}$$

Here,  $y(\tilde{x})$  means  $yx_1x_2\dots x_m$  and  $\tilde{y}(\tilde{x})$  means  $(y_1\tilde{x}, y_2\tilde{x}, \dots, y_n\tilde{x})$  when  $\tilde{x} = (x_1, \dots, x_m)$  and  $\tilde{y} = (y_1, \dots, y_n)$ .

Let us say that a formula  $A$  is realizable if there exists a sequence of  $PV^\omega$ -terms  $\vec{t}$  with free variables among those of  $A$  such that  $IPV^\omega \vdash \vec{t} \textcircled{R} A$ .

Induction on derivations shows that whenever  $IPV^\omega \vdash A$ , we have  $A$  is realizable. The key observation is that an instance of NP-induction can be realised using bounded recursion on notation. Cook and Urquhart add  $A \rightarrow B$  as a conjunct to  $\tilde{x} \textcircled{R} A \rightarrow B$ . This ensures that the converse also holds, i.e., realizable formulas are provable. We do not want this property since in Proposition 3.5 below we do not know whether a formula that we need to realise is provable.

3.3. Interpretation of  $IPV^\omega + EM$  in  $IPV^\omega$

The following is a straightforward application of realisability and appears as Corollary 8.17 in Cook and Urquhart (1993).

**Proposition 3.1.** If  $IPV^\omega \vdash \exists y: o.s(\tilde{x}, y)=0$ , there exists a term  $t(\tilde{x}) : o$  such that  $IPV^\omega \vdash s(\tilde{x}, t(\tilde{x})) = 0$ .

The central result of Cook and Urquhart (1993) is that a similar result holds for  $IPV^\omega + EM$  and hence for Buss' bounded arithmetic  $S_2^1$ .

**Theorem 3.2.** Let  $A(\tilde{x}) \stackrel{\text{def}}{=} \exists y: o.s(\tilde{x}, y)=0$  be a  $\Sigma_1^0$  formula. If  $IPV^\omega + EM \vdash A(\tilde{x})$ , we can find a  $PV^\omega$  term  $t(\tilde{x})$  such that  $IPV^\omega \vdash s(\tilde{x}, t(\tilde{x})) = 0$ .

Cook and Urquhart prove this by first using the negative translation to interpret  $IPV^\omega + EM$  in  $IPV^\omega + MP$  (where  $MP$  stands for Markov's principle), and then interpreting  $IPV^\omega + MP$  in  $IPV^\omega$  using a variant of the Dialectica interpretation.

We will now give an interpretation of  $IPV^\omega + MP$  in  $IPV^\omega$  using our method.

To every formula  $A$  of  $IPV^\omega$  and fresh variable  $f : o \rightarrow o$ , we can associate a formula  $f \Vdash A$  of  $IPV^\omega$  by the clauses given above in Section 2.2. The proof of Lemma 2.4 then goes through without changes and we can conclude that if  $A$  is a (not necessarily bounded)  $\Sigma_1^0$ -formula, then  $f \Vdash \neg\neg A$  is equivalent to  $A \vee T(f)$ , where, as before,  $T(f) \stackrel{\text{def}}{=} \exists x.f(x)=0$ .

In order to establish the analogue of Lemma 2.6, i.e., that  $f \Vdash A$  for all instances of NP-induction  $A$ , we need to strengthen the induction scheme of  $IPV^\omega$  slightly.

Let  $IPV_+^\omega$  stand for  $IPV^\omega$  extended by all formulas  $PIND(A(x) \vee B)$  where  $A$  is bounded  $\Sigma_1^0$  and  $B$  is an arbitrary formula not containing  $x$ . This extension was introduced by Buss in Buss (1990). It is shown there that it is complete for a certain class of Kripke-models.

The following is a direct adaptation of Lemma 2.6.

**Lemma 3.3.** For every bounded  $\Sigma_1^0$  formula  $A$  the formula  $f \Vdash PIND(A)$  is provable in  $IPV_+^\omega$ .

This gives a proof of the following theorem.

**Theorem 3.4.** If  $A$  is provable in  $IPV^\omega + MP$ , then  $f \Vdash A$  is provable in  $IPV_+^\omega$ .

In order to deduce the main result that Skolem functions of  $IPV^\omega + MP$  are *PTIME*, we must extend our polynomial realisability to  $IPV_+^\omega$ .

**Proposition 3.5.** Let  $A(x) \stackrel{\text{def}}{=} \exists y \leq t(x).s(x, y)$  and  $B$  be arbitrary so that  $x$  does not occur in  $B$ . The formula  $PIND(A(x) \vee B)$  is realizable.

*Proof.* The formula  $PIND(A(x) \vee B)$  is equivalent to

$$(A(0) \vee B) \rightarrow (\forall x. \forall z \leq t(\lfloor \frac{x}{2} \rfloor).s(\lfloor \frac{x}{2} \rfloor, z)=0 \rightarrow (\exists y \leq t(x).s(x, y)=0) \vee B) \rightarrow \forall x. A(x) \vee B.$$

To realise the latter formula, assume variables  $u, g_0, \vec{g}_1, v, h_0, \vec{h}_1$  such that (in the sense of  $IPV^\omega$ -assumptions)

$$\begin{aligned} u=0 &\rightarrow g_0 \leq t(0) \wedge s(0, g_0)=0 \\ u=1 &\rightarrow (\vec{g}_1 \text{ @ } B) \\ z \leq t(\lfloor \frac{x}{2} \rfloor) \wedge s(\lfloor \frac{x}{2} \rfloor, z)=0 \wedge v(x, z)=0 &\rightarrow h_0(x, z) \leq t(x) \wedge s(x, h_0(x, z))=0 \\ z \leq t(\lfloor \frac{x}{2} \rfloor) \wedge s(\lfloor \frac{x}{2} \rfloor, z)=0 \wedge v(x, z)=1 &\rightarrow (\vec{h}_1(x, z) \text{ @ } B), \end{aligned}$$

that is to say, we assume that

$$\begin{aligned} u, g_0, \vec{g}_1 &\text{ @ } A(0) \vee B \\ v, h_0, \vec{h}_1 &\text{ @ } \forall x. \forall z \leq t(\lfloor \frac{x}{2} \rfloor).s(\lfloor \frac{x}{2} \rfloor, z)=0 \rightarrow A(x) \vee B. \end{aligned}$$

We must construct functions  $w, f_0, \vec{f}_1$  such that (in  $IPV^\omega$  under the above assumptions)

$$\begin{aligned} w(x)=0 &\rightarrow f_0(x) \leq t(x) \wedge s(x, f_0(x))=0 \\ w(x)=1 &\rightarrow (\vec{f}_1(x) \text{ @ } B). \end{aligned}$$

That is,  $w, f_0, \vec{f}_1 \text{ @ } A(x) \vee B$ .

We will use the  $PV^\omega$ -recursor to define a function  $F : o \rightarrow o$  meeting the following specification:

— if  $F(x)$  is even, then

$$\lfloor F(x)/2 \rfloor \text{ @ } A(x)$$

— if  $F(x)$  is odd, then

$$u=1 \vee (y_0 \leq t(\lfloor \frac{x_0}{2} \rfloor) \wedge v(x_0, y_0)=1 \wedge s(\lfloor \frac{x_0}{2} \rfloor, y_0)=0)$$



where

$$\begin{aligned} x_0 &= \lfloor F(x)/2 \rfloor \\ y_0 &= \lfloor F(\lfloor x_0/2 \rfloor)/2 \rfloor \end{aligned}$$

Notice that in the second case  $\vec{h}_1(x_0, y_0) \textcircled{R} B$  by the assumption on  $\vec{h}_1$ . The reason for not making this consequence part of the specification is that it need not be a  $\Sigma_b^1$ -formula, and hence cannot necessarily be established by NP-induction. The official specification of  $F$ , however, is equivalent to an atomic formula.

Once we have such a function  $F$  we can define the desired functions  $w, f_0, \vec{f}_1$  by

$$\begin{aligned} w(x) &= \text{Parity}(F(x)) \\ f_0(x) &= \lfloor \frac{F(x)}{2} \rfloor \\ \vec{f}_1(x) &= \text{Cond}(u, \vec{g}_1, \vec{h}_1(x_0, y_0)) \\ &\text{where } x_0 = \lfloor F(x)/2 \rfloor \text{ and } y_0 = \lfloor F(\lfloor x_0/2 \rfloor)/2 \rfloor. \end{aligned}$$

The following is a sugared definition of  $F$  by recursion on notation:

$$\begin{aligned} F(0) &= \text{if } u=0 \text{ then } s_0(g_0) \text{ else } s_1(0) \\ F(x) &= \text{if } \text{Parity}(F(\lfloor x/2 \rfloor))=0 \\ &\text{then let } z = \lfloor F(\lfloor x/2 \rfloor)/2 \rfloor \text{ in} \\ &\text{if } v(x, z)=0 \\ &\text{then } s_0(h_0(x, z)) \\ &\text{else } s_1(x) \\ &\text{else } F(\lfloor x/2 \rfloor) \end{aligned}$$

where the second case applies when  $x > 0$ .

Induction on  $x$  (using the defining property of  $h_0$ ) now shows that if  $F(x)$  is even, then  $F(x) \leq s_0(t(x)) \wedge s(x, \lfloor \frac{F(x)}{2} \rfloor) = 0$ . If  $F(x)$  is odd, then, clearly,  $F(x) \leq s_1(x)$ . Therefore, the above can be transformed into a legal  $PV^\omega$ -definition with  $k(x) := \max(s_0(t(x)), s_1(x))$ .

By formalised NP-induction on  $x$ , it is now possible to show that the resulting function  $F$  satisfies its defining equations (because the bound  $k$  is ‘valid’) and the above specification. The required properties of the derived functions  $w, f_0, \vec{f}_1$  are then direct.  $\square$

Putting things together yields a proof of Theorem 3.2.

#### 4. Conclusions

We have presented a new method for establishing conservativity of excluded middle with respect to  $\Pi_2^0$ -formulas in situations where the usual negative translation is not applicable. The method thus provides an alternative to the techniques using cut elimination or *Dialectica*-interpretation that have been used before. We believe that the new method is simpler than these previous methods and that it should give rise to a more efficient program extraction procedure than those. Of course, some concrete examples would have to be carried out to substantiate this claim. Independently, Avigad has applied this method to a system of polynomial strength (Avigad 1998) and to the system with  $\Sigma_1^1$  axiom of choice. Avigad’s result on bounded arithmetic is slightly weaker than ours since it only gives conservativity for bounded  $\Pi_2^0$ -formulas. Accordingly, the detour via  $IPV_+^\omega$  is not needed in his proof.

This method has also been used in the work of Burr (Burr 1998) on fragments of Heyting arithmetic.

Some extensions of the proof for  $PRA^\omega$  are possible. The negative translation provides us with an interpretation of  $PRA^\omega + EM + \Sigma_1^0-AC$  in  $PRA^\omega + EM + \Sigma_1^0-AC$ . It is readily seen that  $\Sigma_1^0-AC$  is valid under our Kripke semantics, which is why we can interpret the latter system in  $PRA^\omega + \Sigma_1^0-AC$ , which is  $\Pi_2^0$ -conservative over  $PRA^\omega$  by realisability.

Let  $EL'_1$  be the system  $EL_1$  of Troelstra (1974) with induction restricted to  $\Sigma_1^0$  formula. We can add the fan theorem (*FAN*) to  $EL'_1 + AC$ . Then we get an interpretation of  $EL'_1 + EM + \Sigma_1^0-AC + WKL$  (weak König's lemma) in  $EL'_1 + FAN + AC$ . The usual method of elimination of choice sequence shows that this system is conservative over  $PRA$  for  $\Pi_2^0$  sentences, reading Troelstra (1974) with  $EL'_1$  instead of  $EL_1$ .

These results may be interesting for comparing some recent intuitionistic proofs of theorems such as Hahn–Banach, Heine–Borel (Cederquist 1997), the existence of prime ideals in their localic formulation using formal topology and their corresponding proofs in reverse mathematics where these theorems are proved classically in systems similar to  $EL'_1 + EM + WKL + \Sigma_1^0-AC$ . For an extension of these results to  $PRA^\omega$  and sharper formulations, using a functional interpretation, see Kohlenbach (1992). Whether our method gives an alternative proof of his result remains unexplored.

## References

- Aczel, P., Simmons, H. and Wainer, S. (eds.) (1992) *Proof Theory*, Cambridge University Press.
- Avigad, J. (1998) Interpreting classical theories in constructive ones. Preprint, available from <http://macduff.andrew.cmu.edu/avigad>.
- Burr, W. (1998) Fragments of Heyting-Arithmetic (to appear in *Journal of Symbolic Logic*).
- Buchholz, W., Feferman, S., Pohlers, W. and Sieg, W. (1981) Iterated Inductive Definitions and Subsystems of Analysis: Recent Proof-Theoretical Studies. *Springer-Verlag Lecture Notes in Mathematics* **897**.
- Buchholz, W. (1977) Some proof-theoretical results in the theories  $ID_v^c$ ,  $ID_v^i$ . Preprint, March 1977.
- Buss, S. R. (1990) On model theory for intuitionistic bounded arithmetic with applications to independence results. In: Buss, S. R. and Scott, P. J. (eds.) *Feasible Mathematics*, Birkhäuser 27–47.
- Buss, S. R. (1986) *Bounded Arithmetic*, Bibliopolis.
- Cederquist, J. (1997) *Formal Topology in Type Theory*. Ph. D. thesis, Chalmers University.
- A. Cobham (1965) The intrinsic computational difficulty of functions. In: Bar-Hillel, Y. (ed.) *Logic, Methodology, and Philosophy of Science II*, Springer Verlag 24–30.
- Cook, S. and Urquhart, A. (1993) Functional interpretations of feasibly constructive arithmetic. *Annals of Pure and Applied Logic* **63** 103–200.
- van Heijenoort, J. (ed.) (1971) *From Frege to Gödel*, Harvard University Press.
- Hilbert, D. (1971) The foundations of mathematics. In: van Heijenoort, J. (ed.) *From Frege to Gödel*, Harvard University Press 465–479.
- Jervell, H. R. (1998) A course in proof theory. (Manuscript, currently (April 1998) available at [www.uio.no/~herman/ptheory.ps](http://www.uio.no/~herman/ptheory.ps).)
- Johannsen, J. (1996) *Schwache Fragmente der Arithmetik und Schwellertschaltkreise beschränkter Tiefe*, Ph. D. thesis, Universität Erlangen.
- Johnstone, P. (1982) *Stone Spaces*, Cambridge University Press.

- Kohlenbach, U. (1992) Effective Bounds from Ineffective Proofs in Analysis: an Application of Functional Interpretation and Majorization. *Journal of Symbolic Logic* 1239–1273.
- Parsons, C. (1970) On a number theoretic choice schema and its relation to induction. In: Kino, Myhill and Vesley (eds.) *Intuitionism and Proof Theory*, North-Holland 459–473.
- Parsons, C. (1972) On  $n$ -quantifier induction. *Journal of Symbolic Logic* 37 466–482.
- Pohlers, W. (1992) A short course in ordinal analysis. In: Aczel, P., Simmons, H. and Wainer, S. (eds.) *Proof Theory*, Cambridge University Press 27–78.
- Sambin, G. (1996) A new and elementary method to represent every complete Boolean algebra. In: Ursini, A. and Agliano, P. (eds.) *Logic and algebra* (Pontignano, 1994). *Lecture Notes in Pure and Appl. Math.*, Dekker, New York 180 655–665.
- Tait, W. (1968) Normal Derivability in Classical Logic. In: Barwise, J. (ed.) *Springer-Verlag Lecture Notes in Mathematics* 72 204–236.
- Troelstra, A. (1974) Note on the Fan Theorem. *Journal of Symbolic Logic* 39 584–596.
- Troelstra, A., and van Dalen, D. (1988) *Constructivism in Mathematics, vol. II*, North-Holland.
- Wainer, S. and Wallen, L. (1992) Basic proof theory. In: Aczel, P., Simmons, H. and Wainer, S. (eds.) *Proof Theory*, Cambridge University Press 1–26.