

# Gödel's functional interpretation and higher-order learning

Thomas Powell

University of Innsbruck

MATHEMATICS FOR COMPUTATION

Niederaltaich

Germany

12 May 2016

# WHAT IS THE COMPUTATIONAL MEANING OF CLASSICAL REASONING?

## WHAT IS THE COMPUTATIONAL MEANING OF CLASSICAL REASONING?

One could consider the ‘finitized’ version of classical theorems:

$A$  : there exists an ideal object  $x$ .

$A_{\text{fin}}$  : there exist finite approximations to  $x$  of arbitrary high quality.

## WHAT IS THE COMPUTATIONAL MEANING OF CLASSICAL REASONING?

One could consider the ‘finitized’ version of classical theorems:

$A$  : there exists an ideal object  $x$ .

$A_{\text{fin}}$  : there exist finite approximations to  $x$  of arbitrary high quality.

Over classical logic,  $A \leftrightarrow A_{\text{fin}}$ .

While ideal objects cannot be effectively constructed, finite approximations to them can.

This talk is about *algorithms* which compute such approximations.

What is the computational meaning of a  $\Pi_3$ -theorem?

$$A \equiv \forall x^X \exists y^Y \forall z^Z P(x, y, z)$$

What is the computational meaning of a  $\Pi_3$ -theorem?

$$A := \forall x^X \exists y^Y \forall z^Z P(x, y, z)$$

In general we cannot hope to produce a direct computable witness for  $\exists y$ . But suppose we double negate and Skolemize:

What is the computational meaning of a  $\Pi_3$ -theorem?

$$A \equiv \forall x^X \exists y^Y \forall z^Z P(x, y, z)$$

In general we cannot hope to produce a direct computable witness for  $\exists y$ . But suppose we double negate and Skolemize:

$$\begin{aligned} \neg A &\leftrightarrow \exists x \forall y \exists z \neg P(x, y, z) \\ &\leftrightarrow \exists x, \xi^{Y \rightarrow Z} \forall y \neg P(x, y, \xi(y)) \end{aligned}$$

What is the computational meaning of a  $\Pi_3$ -theorem?

$$A \equiv \forall x^X \exists y^Y \forall z^Z P(x, y, z)$$

In general we cannot hope to produce a direct computable witness for  $\exists y$ . But suppose we double negate and Skolemize:

$$\begin{aligned} \neg A &\leftrightarrow \exists x \forall y \exists z \neg P(x, y, z) \\ &\leftrightarrow \exists x, \xi^{Y \rightarrow Z} \forall y \neg P(x, y, \xi(y)) \\ \neg \neg A &\leftrightarrow \forall x, \xi \exists y \neg \neg P(x, y, \xi(y)) \\ &\leftrightarrow \forall x, \xi \exists y P(x, y, \xi(y)) \end{aligned}$$



What is the computational meaning of a  $\Pi_3$ -theorem?

$$A := \forall x^X \exists y^Y \forall z^Z P(x, y, z)$$

In general we cannot hope to produce a direct computable witness for  $\exists y$ . But suppose we double negate and Skolemize:

$$\begin{aligned} \neg A &\leftrightarrow \exists x \forall y \exists z \neg P(x, y, z) \\ &\leftrightarrow \exists x, \xi^{Y \rightarrow Z} \forall y \neg P(x, y, \xi(y)) \\ \neg \neg A &\leftrightarrow \forall x, \xi \exists y \neg \neg P(x, y, \xi(y)) \\ &\leftrightarrow \forall x, \xi \exists y P(x, y, \xi(y)) \end{aligned}$$

Over classical logic

$$\underbrace{\forall x \exists y \forall z P(x, y, z)}_{y \text{ ideal (for all } z)} \leftrightarrow \underbrace{\forall x, \xi \exists y P(x, y, \xi(y))}_{y \text{ approximate relative to } \xi}$$

but we can realize the R.H.S.

Let  $f: \mathbb{N} \rightarrow \mathbb{N}$  be an arbitrary function.

INFINITARY. For any  $x \in \mathbb{N}$  there exists some  $y \geq x$  such that

$$\forall z [z \geq x \rightarrow f(y) \leq f(z)].$$

Let  $f: \mathbb{N} \rightarrow \mathbb{N}$  be an arbitrary function.

INFINITARY. For any  $x \in \mathbb{N}$  there exists some  $y \geq x$  such that

$$\forall z [z \geq x \rightarrow f(y) \leq f(z)].$$

FINITARY. For any  $x \in \mathbb{N}$  and  $\xi: \mathbb{N} \rightarrow \mathbb{N}$  there exists some  $y \geq x$  such that

$$\xi(y) \geq y \rightarrow f(y) \leq f(\xi(y)).$$

Let  $f: \mathbb{N} \rightarrow \mathbb{N}$  be an arbitrary function.

INFINITARY. For any  $x \in \mathbb{N}$  there exists some  $y \geq x$  such that

$$\forall z [z \geq x \rightarrow f(y) \leq f(z)].$$

FINITARY. For any  $x \in \mathbb{N}$  and  $\xi: \mathbb{N} \rightarrow \mathbb{N}$  there exists some  $y \geq x$  such that

$$\xi(y) \geq y \rightarrow f(y) \leq f(\xi(y)).$$

We can compute  $y$  by *learning*, as follows:

$$y := \begin{cases} x & \text{if } \xi(x) \geq x \rightarrow f(x) \leq f(\xi(x)) \\ \xi(x) & \text{if } \xi^{(2)}(x) \geq \xi(x) \rightarrow f(\xi(x)) \leq f(\xi^{(2)}(x)) \\ \xi^{(2)}(x) & \text{if } \xi^{(3)}(x) \geq \xi^{(2)}(x) \rightarrow f(\xi^{(2)}(x)) \leq f(\xi^{(3)}(x)) \\ \dots & \dots \end{cases}$$

Terminates since otherwise we'd have  $f(x) > f(\xi(x)) > f(\xi^{(2)}(x)) > \dots$

A more interesting example: Cauchy convergence of a monotone sequence  $(a_i) \in [0, 1]^\omega$

INFINITARY. For any  $x$  there exists  $y$  such that  $i, j \geq y$  implies  $|a_i - a_j| < 2^{-x}$ .

FINITARY (T. TAO). For any  $x$  and  $\xi: \mathbb{N} \rightarrow \mathbb{N}$  there exists  $Y$  such that any  $0 \leq a_0 \leq \dots \leq a_Y \leq 1$  there exists  $y$  with  $0 \leq y < y + \xi(y) \leq Y$  such that  $|a_i - a_j| < 2^{-x}$  for all  $y \leq i, j \leq \xi(y)$ .

Moreover, we can show that  $Y \leq (\lambda i.i + \xi(i))^{(2^x)}(0)$ .

The general problem of constructing a realizer for a finitary formulation of a classical theorem is highly non-trivial. We can appeal to Gödel's functional interpretation (more precisely his ND interpretation):

$$A \rightsquigarrow A^{\neg\neg} \rightsquigarrow \exists x \forall y |A^{\neg\neg}|_y^x$$

which allows us to extract from a proof  $\mathcal{P}$  of  $A$  a term  $t$  satisfying  $\forall y |A^{\neg\neg}|_y^t$

The general problem of constructing a realizer for a finitary formulation of a classical theorem is highly non-trivial. We can appeal to Gödel's functional interpretation (more precisely his ND interpretation):

$$A \rightsquigarrow A^{\neg\neg} \rightsquigarrow \exists x \forall y |A^{\neg\neg}|_y^x$$

which allows us to extract from a proof  $\mathcal{P}$  of  $A$  a term  $t$  satisfying  $\forall y |A^{\neg\neg}|_y^t$

However, a brute force extraction can yield terms which are highly inefficient and difficult to understand.

The general problem of constructing a realizer for a finitary formulation of a classical theorem is highly non-trivial. We can appeal to Gödel's functional interpretation (more precisely his ND interpretation):

$$A \rightsquigarrow A^{\neg\neg} \rightsquigarrow \exists x \forall y |A^{\neg\neg}|_y^x$$

which allows us to extract from a proof  $\mathcal{P}$  of  $A$  a term  $t$  satisfying  $\forall y |A^{\neg\neg}|_y^t$

However, a brute force extraction can yield terms which are highly inefficient and difficult to understand.

Given the wide range of applications of the functional interpretation in modern proof theory, it can be extremely useful to devise refinements of the functional interpretation which help us extract better terms.



The general problem of constructing a realizer for a finitary formulation of a classical theorem is highly non-trivial. We can appeal to Gödel's functional interpretation (more precisely his ND interpretation):

$$A \rightsquigarrow A^{\neg\neg} \rightsquigarrow \exists x \forall y |A^{\neg\neg}|_y^x$$

which allows us to extract from a proof  $\mathcal{P}$  of  $A$  a term  $t$  satisfying  $\forall y |A^{\neg\neg}|_y^t$

However, a brute force extraction can yield terms which are highly inefficient and difficult to understand.

Given the wide range of applications of the functional interpretation in modern proof theory, it can be extremely useful to devise refinements of the functional interpretation which help us extract better terms.

I will sketch one such refinement, which enriches the usual interpreting system of higher type recursors with *learning algorithms* that produce realizing terms which are more efficient and intuitive.

## THE INTUITION BEHIND LEARNING ALGORITHMS

Suppose we have a decidable predicate  $G$  on  $X$ , and we want to find some  $x \in X$  satisfying  $G(x)$ .

## THE INTUITION BEHIND LEARNING ALGORITHMS

Suppose we have a decidable predicate  $G$  on  $X$ , and we want to find some  $x \in X$  satisfying  $G(x)$ .

Let's some arbitrary  $x_0$ . If  $G(x_0)$  holds then we're done. On the other hand, if  $\neg G(x_0)$  and we fail we often learn a useful piece of constructive information  $\xi(x_0)$ .

## THE INTUITION BEHIND LEARNING ALGORITHMS

Suppose we have a decidable predicate  $G$  on  $X$ , and we want to find some  $x \in X$  satisfying  $G(x)$ .

Let's some arbitrary  $x_0$ . If  $G(x_0)$  holds then we're done. On the other hand, if  $\neg G(x_0)$  and we fail we often learn a useful piece of constructive information  $\xi(x_0)$ .

We can then use this to update our original guess with a better one  $x_1 := x_0 \oplus \xi(x_0)$ , and continue:

## THE INTUITION BEHIND LEARNING ALGORITHMS

Suppose we have a decidable predicate  $G$  on  $X$ , and we want to find some  $x \in X$  satisfying  $G(x)$ .

Let's some arbitrary  $x_0$ . If  $G(x_0)$  holds then we're done. On the other hand, if  $\neg G(x_0)$  and we fail we often learn a useful piece of constructive information  $\xi(x_0)$ .

We can then use this to update our original guess with a better one  $x_1 := x_0 \oplus \xi(x_0)$ , and continue:

$$x := \begin{cases} x_0 & \text{if } G(x_0) \\ x_1 := x_0 \oplus \xi(x_0) & \text{if } G(x_1) \\ x_2 := x_1 \oplus \xi(x_1) & \text{if } G(x_2) \\ \dots & \dots \end{cases}$$

The idea is that we eventually reach some  $x_k$  satisfying  $G(x_k)$ .

## LEARNING ALGORITHMS - A FORMAL DEFINITION

A learning algorithm of type  $X, L$  is a tuple  $\mathcal{L} = (\mathbf{G}, \xi, \oplus)$  where

- $\mathbf{G} : X \rightarrow \mathbb{B}$  is a decidable predicate which tests whether an element  $x \in X$  is ‘good’;
- $\xi : X \rightarrow L$  and  $\oplus : X \times L \rightarrow X$  are responsible for learning, and will be used to map bad objects  $x \in X$  to improvements  $x \oplus \xi(x)$ ;

## LEARNING ALGORITHMS - A FORMAL DEFINITION

A learning algorithm of type  $X, L$  is a tuple  $\mathcal{L} = (\mathbf{G}, \xi, \oplus)$  where

- $\mathbf{G} : X \rightarrow \mathbb{B}$  is a decidable predicate which tests whether an element  $x \in X$  is ‘good’;
- $\xi : X \rightarrow L$  and  $\oplus : X \times L \rightarrow X$  are responsible for learning, and will be used to map bad objects  $x \in X$  to improvements  $x \oplus \xi(x)$ ;

The learning procedure  $\mathcal{L}[x]$  starting at  $x \in X$  is a sequence  $(x_i) \in X^{\mathbb{N}}$  defined by

$$x_0 := x \quad \text{and} \quad x_{i+1} := \begin{cases} x_i & \text{if } \mathbf{G}(x_i) \\ x_i \oplus \xi(x_i) & \text{if } \neg \mathbf{G}(x_i) \end{cases}$$

## LEARNING ALGORITHMS - A FORMAL DEFINITION

A learning algorithm of type  $X, L$  is a tuple  $\mathcal{L} = (\mathbf{G}, \xi, \oplus)$  where

- $\mathbf{G} : X \rightarrow \mathbb{B}$  is a decidable predicate which tests whether an element  $x \in X$  is ‘good’;
- $\xi : X \rightarrow L$  and  $\oplus : X \times L \rightarrow X$  are responsible for learning, and will be used to map bad objects  $x \in X$  to improvements  $x \oplus \xi(x)$ ;

The learning procedure  $\mathcal{L}[x]$  starting at  $x \in X$  is a sequence  $(x_i) \in X^{\mathbb{N}}$  defined by

$$x_0 := x \quad \text{and} \quad x_{i+1} := \begin{cases} x_i & \text{if } \mathbf{G}(x_i) \\ x_i \oplus \xi(x_i) & \text{if } \neg \mathbf{G}(x_i) \end{cases}$$

The limit of  $\mathcal{L}[x]$  is defined as

$$\lim \mathcal{L}[x] := x_k$$

where  $x_k$  is the least point satisfying  $\mathbf{G}(x_k)$  (whenever it exists).



VAGUE IDEA. Suppose that  $\exists xA(x)$  is a classical theorem, and

$$\forall \xi \exists x A'_\xi(x, \xi)$$

a finitization of  $A$ . Then  $x$  can be computed in the limit of a learning procedure parametrised by  $\xi$ :

$$x := F(\lim \mathcal{L}_\xi[x_0])$$

for some  $x_0$ .

VAGUE IDEA. Suppose that  $\exists x A(x)$  is a classical theorem, and

$$\forall \xi \exists x A'_\xi(x, \xi)$$

a finitization of  $A$ . Then  $x$  can be computed in the limit of a learning procedure parametrised by  $\xi$ :

$$x := F(\lim \mathcal{L}_\xi[x_0])$$

for some  $x_0$ .

This idea is made more precise in (P. 2016), where a collection of concrete results of this kind are given, relating Gödel's functional interpretation of induction and comprehension principles to learning procedures.

In the remainder of the talk, however, I will not go into any detail, and will simply give some illustrations.

EXAMPLE 1: The quantifier-free minimum principle

$$\text{QFMin} : \exists x P(x) \rightarrow \exists y (P(y) \wedge \forall z \prec y \neg P(z)).$$

EXAMPLE 1: The quantifier-free minimum principle

$$\text{QFMin} : \exists x P(x) \rightarrow \exists y (P(y) \wedge \forall z \prec y \neg P(z)).$$

This has an ND interpretation given by

$$(*) \quad \forall x, \xi \exists y (P(x) \rightarrow P(y) \wedge (\xi(y) \prec y \rightarrow \neg P(\xi(y))))$$

*“For all  $x, \xi$  there exists some  $y$  such that whenever  $P(x)$  holds then  $P(y)$  holds and  $y$  is approximately minimal with respect to  $\xi(y)$ ”*

EXAMPLE 1: The quantifier-free minimum principle

$$\text{QFMin} : \exists x P(x) \rightarrow \exists y (P(y) \wedge \forall z \prec y \neg P(z)).$$

This has an ND interpretation given by

$$(*) \quad \forall x, \xi \exists y (P(x) \rightarrow P(y) \wedge (\xi(y) \prec y \rightarrow \neg P(\xi(y))))$$

*“For all  $x$ ,  $\xi$  there exists some  $y$  such that whenever  $P(x)$  holds then  $P(y)$  holds and  $y$  is approximately minimal with respect to  $\xi(y)$ ”*

We can compute  $y$  in  $x$  and  $\xi$  using the following idea

$$y := \begin{cases} x & \text{if } \xi(x) \prec x \rightarrow \neg P(\xi(x)) \\ \xi(x) & \text{if } \xi^{(2)}(x) \prec \xi(x) \rightarrow \neg P(\xi^{(2)}(x)) \\ \xi^{(2)}(x) & \text{if } \dots \\ \dots & \dots \end{cases}$$

THEOREM. Define  $\mathcal{L}_\xi := (\mathbf{G}_\xi, \xi, \pi_2)$  where

$$\mathbf{G}_\xi(x) \leftrightarrow [\xi(x) \prec x \rightarrow \neg P(\xi(x))].$$

Then the ND interpretation of QFMin, given by

$$\forall x, \xi \exists y (P(x) \rightarrow P(y) \wedge (\xi(y) \prec y \rightarrow \neg P(\xi(y))))$$

is realized by

$$\lambda x, \xi . \lim \mathcal{L}_\xi[x].$$

REMARK. A general result dealing with well-founded induction for arbitrary formulas and relations  $\prec$  is given in (P. 2016), inspired by (Schwichtenberg 2008).

EXAMPLE 2, FOLLOWING (SCHWICHTENBERG 2008). For any two natural numbers  $a, b > 0$  there exist integers  $m, n$  such that  $am + bn \mid a, b$ .

Classical proof. Use a variant of QFMin relative to the ordering  $(x, y) \prec (x', y') := ax + by < ax' + by'$ .

EXAMPLE 2, FOLLOWING (SCHWICHTENBERG 2008). For any two natural numbers  $a, b > 0$  there exist integers  $m, n$  such that  $am + bn \mid a, b$ .

Classical proof. Use a variant of QFMin relative to the ordering  $(x, y) \prec (x', y') := ax + by < ax' + by'$ .

A program for computing  $m, n$  in  $a, b$  can be extracted, namely a learning procedure of type  $(\mathbb{N}^{(2)})^*, \mathbb{N}^{(2)}$  given by

$$(m, n) := \text{tail}(\lim \mathcal{L}_{a,b}[\langle e_0, e_1 \rangle])$$

where  $\mathcal{L}_{a,b} = (\mathbf{G}_{a,b}, \xi_{a,b}, ::)$  for

- $\mathbf{G}_{a,b}(s) \leftrightarrow \text{rem}(s_{l-2} \cdot (a, b), s_{l-1} \cdot (a, b)) = 0$
- $\xi_{a,b}(s) \leftrightarrow s_{l-2} - \text{quot}(s_{l-2} \cdot (a, b), s_{l-1} \cdot (a, b))s_{l-1}$
- $s :: x := \langle s_0, \dots, s_{l-1}, x \rangle$ .



EXAMPLE 2, FOLLOWING (SCHWICHTENBERG 2008). For any two natural numbers  $a, b > 0$  there exist integers  $m, n$  such that  $am + bn \mid a, b$ .

Classical proof. Use a variant of QFMin relative to the ordering  $(x, y) \prec (x', y') := ax + by < ax' + by'$ .

A program for computing  $m, n$  in  $a, b$  can be extracted, namely a learning procedure of type  $(\mathbb{N}^{(2)})^*, \mathbb{N}^{(2)}$  given by

$$(m, n) := \text{tail}(\lim \mathcal{L}_{a,b}[\langle e_0, e_1 \rangle])$$

where  $\mathcal{L}_{a,b} = (\mathbf{G}_{a,b}, \xi_{a,b}, ::)$  for

- $\mathbf{G}_{a,b}(s) \leftrightarrow \text{rem}(s_{l-2} \cdot (a, b), s_{l-1} \cdot (a, b)) = 0$
- $\xi_{a,b}(s) \leftrightarrow s_{l-2} - \text{quot}(s_{l-2} \cdot (a, b), s_{l-1} \cdot (a, b))s_{l-1}$
- $s :: x := \langle s_0, \dots, s_{l-1}, x \rangle$ .

This is just the Euclidean algorithm!

## META-EXAMPLE: THE DOUBLE NEGATION SHIFT

$$\text{DNS : } \forall n \neg \neg \exists x^X \forall y P_n(x, y) \rightarrow \neg \neg \forall n \exists x \forall y P_n(x, y)$$

## META-EXAMPLE: THE DOUBLE NEGATION SHIFT

$$\text{DNS} : \forall n \neg \neg \exists x^X \forall y P_n(x, y) \rightarrow \neg \neg \forall n \exists x \forall y P_n(x, y)$$

This has a (partial) functional interpretation given by

$$\forall n, \xi \exists x P_n(x, \xi(x)) \rightarrow \forall \omega, \varphi \exists \alpha^{X^{\mathbb{N}}} P_{\omega\alpha}(\alpha(\omega\alpha), \varphi\alpha)$$

*“If, for each  $n$ ,  $P_n$  is approximately witnessed by some  $x \in X$  relative to  $\xi$ , then we can produce a ‘global’ witness  $\alpha \in X^{\mathbb{N}}$  for the conclusion which is approximately correct relative to  $\varphi\alpha$  at point  $\omega\alpha$ ”*

## META-EXAMPLE: THE DOUBLE NEGATION SHIFT

$$\text{DNS} : \forall n \neg \neg \exists x^X \forall y P_n(x, y) \rightarrow \neg \neg \forall n \exists x \forall y P_n(x, y)$$

This has a (partial) functional interpretation given by

$$\underbrace{\forall n, \xi \exists x P_n(x, \xi(x))}_{(\mathcal{L}_{n, \xi})_{n < \infty}} \rightarrow \underbrace{\forall \omega, \varphi \exists \alpha^{X^{\mathbb{N}}} P_{\omega \alpha}(\alpha(\omega \alpha), \varphi \alpha)}_{\mathcal{L}_{\infty, (\omega, \varphi)}}$$

*“If, for each  $n$ ,  $P_n$  is approximately witnessed by some  $x \in X$  relative to  $\xi$ , then we can produce a ‘global’ witness  $\alpha \in X^{\mathbb{N}}$  for the conclusion which is approximately correct relative to  $\varphi \alpha$  at point  $\omega \alpha$ ”*

Is there some sort of formal construction from a collection of ‘pointwise’ learning algorithms to a ‘global’ learning algorithm:

$$(\mathcal{L}_{n, \xi})_{n < \infty} \mapsto \mathcal{L}_{\infty, (\omega, \phi)}?$$

## META-EXAMPLE: THE DOUBLE NEGATION SHIFT

$$\text{DNS} : \forall n \neg \neg \exists x^X \forall y P_n(x, y) \rightarrow \neg \neg \forall n \exists x \forall y P_n(x, y)$$

This has a (partial) functional interpretation given by

$$\underbrace{\forall n, \xi \exists x P_n(x, \xi(x))}_{(\mathcal{L}_{n, \xi})_{n < \infty}} \rightarrow \underbrace{\forall \omega, \varphi \exists \alpha^{X^{\mathbb{N}}} P_{\omega \alpha}(\alpha(\omega \alpha), \varphi \alpha)}_{\mathcal{L}_{\infty, (\omega, \varphi)}}$$

*“If, for each  $n$ ,  $P_n$  is approximately witnessed by some  $x \in X$  relative to  $\xi$ , then we can produce a ‘global’ witness  $\alpha \in X^{\mathbb{N}}$  for the conclusion which is approximately correct relative to  $\varphi \alpha$  at point  $\omega \alpha$ ”*

Is there some sort of formal construction from a collection of ‘pointwise’ learning algorithms to a ‘global’ learning algorithm:

$$(\mathcal{L}_{n, \xi})_{n < \infty} \mapsto \mathcal{L}_{\infty, (\omega, \phi)}?$$

YES! Details in (P. 2016), I will just give an illustration here.

### EXAMPLE 3: LAW OF EXCLUDED MIDDLE FOR $\Sigma_1^0$ -FORMULAS

$$\begin{aligned} & \forall n \exists b^{\mathbb{B}} (\exists x Q_n(x) \vee_b \forall y \neg Q_n(y)) \\ \rightsquigarrow & \forall n \exists b, x \forall y (Q_n(x) \vee \neg Q_n(y)) \end{aligned}$$

### EXAMPLE 3: LAW OF EXCLUDED MIDDLE FOR $\Sigma_1^0$ -FORMULAS

$$\begin{aligned} & \forall n \exists b^{\mathbb{B}} (\exists x Q_n(x) \vee_b \forall y \neg Q_n(y)) \\ \rightsquigarrow & \forall n \exists b, x \forall y (Q_n(x) \vee \neg Q_n(y)) \end{aligned}$$

This has a functional interpretation given by

$$\forall n, \xi \exists b, x (Q_n(x) \vee \neg Q_n(\xi(b, x)))$$

### EXAMPLE 3: LAW OF EXCLUDED MIDDLE FOR $\Sigma_1^0$ -FORMULAS

$$\begin{aligned} & \forall n \exists b^{\mathbb{B}} (\exists x Q_n(x) \vee_b \forall y \neg Q_n(y)) \\ \rightsquigarrow & \forall n \exists b, x \forall y (Q_n(x) \vee \neg Q_n(y)) \end{aligned}$$

This has a functional interpretation given by

$$\forall n, \xi \exists b, x (Q_n(x) \vee \neg Q_n(\xi(b, x)))$$

which is realized by a learning procedure of length at most two:

$$b, x = \begin{cases} \perp, 0 & \text{if } \neg Q_n(\xi(\perp, 0)) \\ \top, \xi(\perp, 0) & \text{otherwise.} \end{cases}$$



#### EXAMPLE 4: $\Sigma_1^0$ -ARITHMETIC COMPREHENSION

$$\exists \alpha^{\mathbb{N} \rightarrow \mathbb{B} \times \mathbb{N}} \forall n \forall y [Q_n(\alpha(n)_1) \vee_{\alpha(n)_0} \neg Q_n(y)]$$

#### EXAMPLE 4: $\Sigma_1^0$ -ARITHMETIC COMPREHENSION

$$\exists \alpha^{\mathbb{N} \rightarrow \mathbb{B} \times \mathbb{N}} \forall n \forall y [Q_n(\alpha(n)_1) \vee_{\alpha(n)_0} \neg Q_n(y)]$$

i.e.  $\alpha(n)_0 \in \mathbb{B}$  indicates the truth of  $\exists x Q_n(x)$  and  $\alpha(n)_1 \in \mathbb{N}$  provides a witness. This has a functional interpretation given by

$$\forall \omega, \varphi \exists \alpha [Q_{\omega\alpha}(\alpha(\omega\alpha)_1) \vee_{\alpha(\omega\alpha)_0} \neg Q_{\omega\alpha}(\varphi\alpha)]$$

#### EXAMPLE 4: $\Sigma_1^0$ -ARITHMETIC COMPREHENSION

$$\exists \alpha^{\mathbb{N} \rightarrow \mathbb{B} \times \mathbb{N}} \forall n \forall y [Q_n(\alpha(n)_1) \vee_{\alpha(n)_0} \neg Q_n(y)]$$

i.e.  $\alpha(n)_0 \in \mathbb{B}$  indicates the truth of  $\exists x Q_n(x)$  and  $\alpha(n)_1 \in \mathbb{N}$  provides a witness. This has a functional interpretation given by

$$\forall \omega, \varphi \exists \alpha [Q_{\omega\alpha}(\alpha(\omega\alpha)_1) \vee_{\alpha(\omega\alpha)_0} \neg Q_{\omega\alpha}(\varphi\alpha)]$$

which can be realized by a learning procedure on  $\mathbb{N} \rightarrow \mathbb{B} \times \mathbb{N}$  of unbounded length, but of ‘pointwise’ length at most two.

$$\alpha = \begin{cases} \underbrace{[\ ]}_{\alpha_0} := \lambda n. (\perp, 0) & \text{if } \neg Q_{\omega\alpha_0}(\varphi\alpha_0) \\ \underbrace{[\omega\alpha_0 \mapsto (\top, \varphi\alpha_0)]}_{\alpha_1} & \text{if } \alpha_1(\omega\alpha_1)_0 = \perp \rightarrow \neg Q_{\omega\alpha_1}(\varphi\alpha_1) \\ \underbrace{[\omega\alpha_0 \mapsto (\top, \varphi\alpha_0), \omega\alpha_1 \mapsto (\top, \varphi\alpha_1)]}_{\alpha_2} & \text{if } \alpha_2(\omega\alpha_2)_0 = \perp \rightarrow \neg Q_{\omega\alpha_2}(\varphi\alpha_2) \\ \dots & \dots \end{cases}$$

SUMMARY. For the instance of DNS given by

$$\forall n \neg \neg \exists b, x \forall y [Q_n(x) \vee_b \neg Q_n(y)] \rightarrow \neg \neg \forall n \exists b, x \forall y [Q_n(x) \vee_b Q_n(y)]$$

we have  $(\mathcal{L}_{n,\xi})_{n < \infty}$  given by

$$b, x = \begin{cases} \perp, 0 & \text{if } \neg Q_n(\xi(\perp, 0)) \\ \top, \xi(\perp, 0) & \text{otherwise.} \end{cases}$$

and  $\mathcal{L}_{\infty,(\omega,\varphi)}$  given by

$$\alpha = \begin{cases} [] = \lambda n. (\perp, 0) =: \alpha_0 & \text{if } \neg Q_{\omega\alpha_0}(\varphi\alpha_0) \\ [\omega\alpha_0 \mapsto (\top, \varphi\alpha_0)] =: \alpha_1 & \text{if } \alpha_1(\omega\alpha_1)_0 = \perp \rightarrow \neg Q_{\omega\alpha_1}(\varphi\alpha_1) \\ [\omega\alpha_0 \mapsto (\top, \varphi\alpha_0), \omega\alpha_1 \mapsto (\top, \varphi\alpha_1)] & \text{if } \dots \\ \dots & \dots \end{cases}$$

**Historical aside.** Gödel developed his functional interpretation over a period of 30 years, finally publishing it in 1958. The original paper dealt only with Peano arithmetic, and was extended to analysis by Spector in 1962.

Spector showed that the functional interpretation of the double negation shift (and hence the ND interpretation of countable choice) could be realized using bar recursion in all finite types.

However, Spector's paper was left unfinished - in particular, Section 12.1 sketches an *alternative* to bar recursion for a very simple case of DNS, namely a learning algorithm similar to that presented here. Kreisel received this as a letter, writing:

*“The typescript ... consists of sections up to and including 12.1 of the present paper, with about half a page (crossed out) of a projected 12.2. This last half page states that the proof of the Gödel translation of axiom F [the DNS] would use a generalization of Hilbert's substitution method as illustrated in the special case of 12.1. However Spector's notes do not contain any details, so that it is not quite clear how to reconstruct the proof he had in mind.”*

## DIRECTIONS FOR THE FUTURE

- To establish new complexity results on the length of learning procedures extracted from proofs.
- Try to link our approach to other computational interpretations of classical logic based on learning e.g. the Aschieri-Berardi learning realizability (Aschieri-Berardi 2010).
- Formalise everything so that learning procedures can be automatically extracted, and resulting programs decorated with information indicating how they behave.
- Extend the idea to stronger subsystems of mathematics. I conjecture that weak forms of Zorn's lemma of the kind used to prove Kruskal's theorem can be interpreted via a learning procedure which computes approximations to maximal elements in chain complete posets.

## REFERENCES

**Aschieri-Berardi 2010** Interactive learning-based realizability for Heyting arithmetic with EM1. *Logical Methods in Computer Science*, 2010

**Gödel 1958** Über eine bisher noch nicht benützte Erweiterung des finiten Standpunktes. *dialectica*, 1958

**Powell 2016**. Gödel's functional interpretation and the concept of learning. *LICS 2016*

**Schwichtenberg 2008**. Dialectica interpretation of well founded induction. *Math. Log. Quart.* 2008

**Spector 1962**. Provably recursive functionals of analysis: a consistency proof of analysis by an extension of principles in current intuitionistic mathematics. *Recursive Function Theory: Proc. Symposia in Pure Mathematics*. American Mathematical Society, 1962