

# Gödel's functional interpretation and the concept of learning

Thomas Powell

University of Innsbruck

WORKSHOP ON  
EFFICIENT AND NATURAL PROOF SYSTEMS

University of Bath, UK

16 December 2015

Notation.  $\Sigma_2^0$  formulas written as  $P \equiv \exists x \forall y |P|_y^x$ , where  $|P|_y^x$  is always decidable.

Notation.  $\Sigma_2^0$  formulas written as  $P \equiv \exists x \forall y |P|_y^x$ , where  $|P|_y^x$  is always decidable.

A SIMPLE NON-CONSTRUCTIVE THEOREM:  $\exists n \forall m (|P|^m \rightarrow |P|^n)$

“In a pub there is always a person such that if anyone else is drinking, then that person is drinking”

Notation.  $\Sigma_2^0$  formulas written as  $P \equiv \exists x \forall y |P|_y^x$ , where  $|P|_y^x$  is always decidable.

A SIMPLE NON-CONSTRUCTIVE THEOREM:  $\exists n \forall m (|P|^m \rightarrow |P|^n)$

“In a pub there is always a person such that if anyone else is drinking, then that person is drinking”

$$\frac{\frac{\frac{\exists k |P|^k \vee \forall k \neg |P|^k}{|P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)}{\exists k |P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)} \quad \frac{\frac{\forall k \neg |P|^k \rightarrow |P|^m \rightarrow |P|^0}{\forall k \neg |P|^k \rightarrow \forall m (|P|^m \rightarrow |P|^0)}}{\forall k \neg |P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)}}{\exists n \forall m (|P|^m \rightarrow |P|^n)}$$

Notation.  $\Sigma_2^0$  formulas written as  $P \equiv \exists x \forall y |P|_y^x$ , where  $|P|_y^x$  is always decidable.

A SIMPLE NON-CONSTRUCTIVE THEOREM:  $\exists n \forall m (|P|^m \rightarrow |P|^n)$

“In a pub there is always a person such that if anyone else is drinking, then that person is drinking”

$$\frac{\frac{\frac{\exists k |P|^k \vee \forall k \neg |P|^k}{|P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)}{\exists k |P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)} \quad \frac{\frac{\forall k \neg |P|^k \rightarrow |P|^m \rightarrow |P|^0}{\forall k \neg |P|^k \rightarrow \forall m (|P|^m \rightarrow |P|^0)}}{\forall k \neg |P|^k \rightarrow \exists n \forall m (|P|^m \rightarrow |P|^n)}}{\exists n \forall m (|P|^m \rightarrow |P|^n)}$$

There is no *effective* way of realizing  $\exists n$ . So what is the constructive interpretation of the drinkers paradox?

## METHOD I: HILBERT'S $\epsilon$ -CALCULUS

IDEA: Replace quantifiers by ‘magic’  $\epsilon$ -terms:

$$\exists k A(k) \rightsquigarrow A(\epsilon_k A),$$

and quantifier axioms by critical formulas:

$$A(t) \rightarrow A(\epsilon_k A).$$

1. *Translation.* Convert proofs in predicate logic to proofs in epsilon calculus. Instances of quantifier axioms replaced by critical formulas.

2. *Epsilon elimination.* Suppose we only use a finite set of critical formulas. Interpret all  $\epsilon$ -terms by 0. If we find a mistake i.e.  $A(t) \wedge \neg A(0)$ , ‘learn’ from this mistake and update  $\epsilon_k A \mapsto t$ .

Interpreted proof:

$$\frac{\frac{\frac{\exists k|P|^k \vee \forall k\neg|P|^k}{|P|^k \rightarrow \forall m(|P|^m \rightarrow |P|^k)} \quad \frac{\frac{\neg|P|^m \rightarrow |P|^m \rightarrow |P|^0}{\forall k\neg|P|^k \rightarrow |P|^m \rightarrow |P|^0}}{\forall k\neg|P|^k \rightarrow \forall m(|P|^m \rightarrow |P|^0)}}{\frac{\frac{|P|^k \rightarrow \exists n\forall m(|P|^m \rightarrow |P|^n)}{\exists k|P|^k \rightarrow \exists n\forall m(|P|^m \rightarrow |P|^n)} \quad \frac{\forall k\neg|P|^k \rightarrow \exists n\forall m(|P|^m \rightarrow |P|^n)}}{\exists n\forall m(|P|^m \rightarrow |P|^n)}}$$

Critical formulas:

$\epsilon$ -elimination:

Interpreted proof:

$$\frac{\frac{|P|^{\epsilon_k} \vee \neg|P|^{\epsilon_k} \quad \frac{|P|^{\epsilon_k} \rightarrow \forall m(|P|^m \rightarrow |P|^{\epsilon_k})}{|P|^{\epsilon_k} \rightarrow \exists n \forall m(|P|^m \rightarrow |P|^n)}}{\exists n \forall m(|P|^m \rightarrow |P|^n)} \quad \frac{\frac{\frac{\neg|P|^m \rightarrow |P|^m \rightarrow |P|^0}{\neg|P|^{\epsilon_k} \rightarrow |P|^m \rightarrow |P|^0}}{\neg|P|^{\epsilon_k} \rightarrow \forall m(|P|^m \rightarrow |P|^0)}}{\neg|P|^{\epsilon_k} \rightarrow \exists n \forall m(|P|^m \rightarrow |P|^n)}}$$

Critical formulas:

$$|P|^m \rightarrow |P|^{\epsilon_k}$$

$\epsilon$ -elimination:



Interpreted proof:

$$\frac{\frac{|P|^{\epsilon_k} \vee \neg|P|^{\epsilon_k} \quad \frac{|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m \epsilon_k} \rightarrow |P|^{\epsilon_k}}{|P|^{\epsilon_k} \rightarrow \exists n(|P|^{\epsilon_m n} \rightarrow |P|^n)}}{\exists n(|P|^{\epsilon_m n} \rightarrow |P|^n)} \quad \frac{\frac{\neg|P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}{\neg|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}}{\neg|P|^{\epsilon_k} \rightarrow \exists n(|P|^{\epsilon_m n} \rightarrow |P|^n)}}{\exists n(|P|^{\epsilon_m n} \rightarrow |P|^n)}$$

Critical formulas:

$$|P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_k}$$

$\epsilon$ -elimination:

Interpreted proof:

$$\frac{\frac{|P|^{\epsilon_k} \vee \neg|P|^{\epsilon_k} \quad \frac{|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m \epsilon_k} \rightarrow |P|^{\epsilon_k}}{|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})} \quad \frac{\frac{\neg|P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}{\neg|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}}{\neg|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}}$$

Critical formulas:

$$\begin{aligned} & |P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_k} \\ & (|P|^{\epsilon_m \epsilon_k} \rightarrow |P|^{\epsilon_k}) \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}) \\ & (|P|^{\epsilon_m 0} \rightarrow |P|^0) \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}) \end{aligned}$$

$\epsilon$ -elimination:

Interpreted proof:

$$\frac{\frac{|P|^{\epsilon_k} \vee \neg|P|^{\epsilon_k} \quad \frac{|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m \epsilon_k} \rightarrow |P|^{\epsilon_k}}{|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}} \quad \frac{\frac{\neg|P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}{\neg|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}}{\neg|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}}$$

Critical formulas:

$$\begin{array}{ll} |P|^{\epsilon_m 0} \rightarrow |P|^0 & ? \\ (|P|^{\epsilon_m 0} \rightarrow |P|^0) \rightarrow (|P|^{\epsilon_m 0} \rightarrow |P|^0) & \checkmark \\ (|P|^{\epsilon_m 0} \rightarrow |P|^0) \rightarrow (|P|^{\epsilon_m 0} \rightarrow |P|^0) & \checkmark \end{array}$$

$\epsilon$ -elimination:

- Try  $\epsilon_k = \epsilon_n = 0$ . Works unless  $|P|^{\epsilon_m 0} \wedge \neg|P|^0$ .

Interpreted proof:

$$\frac{\frac{|P|^{\epsilon_k} \vee \neg|P|^{\epsilon_k} \quad \frac{|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m \epsilon_k} \rightarrow |P|^{\epsilon_k}}{|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}} \quad \frac{\frac{\neg|P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}{\neg|P|^{\epsilon_k} \rightarrow |P|^{\epsilon_m 0} \rightarrow |P|^0}}{\neg|P|^{\epsilon_k} \rightarrow (|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n})}}{|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}}$$

Critical formulas:

$$\begin{aligned} & |P|^{\epsilon_m 0} \rightarrow |P|^{\epsilon_m 0} && \checkmark \\ & (|P|^{\epsilon_m \epsilon_m 0} \rightarrow |P|^{\epsilon_m 0}) \rightarrow (|P|^{\epsilon_m \epsilon_m 0} \rightarrow |P|^{\epsilon_m 0}) && \checkmark \\ & (|P|^{\epsilon_m 0} \rightarrow |P|^0) \rightarrow (|P|^{\epsilon_m \epsilon_m 0} \rightarrow |P|^{\epsilon_m 0}) && \checkmark \end{aligned}$$

$\epsilon$ -elimination:

- Try  $\epsilon_k = \epsilon_n = 0$ . Works unless  $|P|^{\epsilon_m 0} \wedge \neg|P|^0$ .
- But now we have a witness for  $\exists k|P|^k$ , so set  $\epsilon_k = \epsilon_m = \epsilon_m 0$ .

FINITARY DRINKER'S PARADOX I: For an arbitrary  $\epsilon$ -term  $\epsilon_m(\cdot)$  there exists some  $\epsilon_n$  satisfying

$$|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}.$$

This can be computed by the algorithm

- Set  $\epsilon_n := 0$ .
- Check  $|P|^{\epsilon_m 0} \rightarrow |P|^0$ . If true, END.
- Else  $\epsilon_n := \epsilon_m 0$ .

FINITARY DRINKER'S PARADOX I: For an arbitrary  $\epsilon$ -term  $\epsilon_m(\cdot)$  there exists some  $\epsilon_n$  satisfying

$$|P|^{\epsilon_m \epsilon_n} \rightarrow |P|^{\epsilon_n}.$$

This can be computed by the algorithm

- Set  $\epsilon_n := 0$ .
- Check  $|P|^{\epsilon_m 0} \rightarrow |P|^0$ . If true, END.
- Else  $\epsilon_n := \epsilon_m 0$ .

The term  $\epsilon_m(\cdot)$  represents the *proof theoretic environment*, a measure of how we might use the drinkers paradox as a lemma. More specifically, exactly when we need the  $\forall$ -axiom

$$\exists n \forall m (|P|^m \rightarrow |P|^n) \rightarrow \exists n (|P|^t \rightarrow |P|^n).$$

## METHOD II: GÖDEL'S FUNCTIONAL (DIALECTICA) INTERPRETATION

IDEA: A two stage translation: *Negative translation* + *Dialectica interpretation*.

1. Eliminate classical reasoning by applying negative translation (can be more flexible here e.g. ignore atomic formulas).
2. Extract realizing terms for *Dialectica* interpretation of this formula. More complex than realizability - need to fully Skolemize implication:

$$\begin{aligned}(A \rightarrow B) &\rightsquigarrow (\exists x \forall y |A|_y^x \rightarrow \exists u \forall v |B|_v^u) \\ &\rightsquigarrow \forall x \exists u \forall v \exists y (|A|_y^x \rightarrow |B|_v^u) \\ &\rightsquigarrow \exists U, Y \forall x, v (|A|_{Yxv}^x \rightarrow |B|_v^{Ux})\end{aligned}$$

Contraction problem: Interpretation of classical reasoning requires us to test atomic formulas and use case definitions.

$$\begin{array}{ccc}
 & [\exists k|P|^k] & [\forall k\neg|P|^k] \\
 & \vdots & \vdots \\
 \exists k|P|^k \vee \forall k\neg|P|^k & \exists n\forall m(|P|^m \rightarrow |P|^n) & \exists n\forall m(|P|^m \rightarrow |P|^n) \\
 \hline
 & \exists n\forall m(|P|^m \rightarrow |P|^n) &
 \end{array}$$



$$\begin{array}{c}
 [\exists k | P|^k] \\
 \vdots \\
 \exists k | P|^k \vee \forall k \neg | P|^k \\
 \neg \neg \exists n \forall m (| P|^m \rightarrow | P|^n) \\
 \hline
 \neg \neg \exists n \forall m (| P|^m \rightarrow | P|^n)
 \end{array}
 \qquad
 \begin{array}{c}
 [\forall k \neg | P|^k] \\
 \vdots \\
 \neg \neg \exists n \forall m (| P|^m \rightarrow | P|^n)
 \end{array}$$

$$\begin{array}{c}
\begin{array}{cc}
[|P|^k] & [\forall k \neg |P|^k] \\
\vdots & \vdots \\
\exists k |P|^k \vee \forall k \neg |P|^k & |P|^{gk} \rightarrow |P|^k \quad \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n)
\end{array} \\
\hline
\neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n)
\end{array}$$

First branch:

$$\begin{aligned}
& \exists k |P|^k \rightarrow \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n) \\
& \rightsquigarrow \exists k |P|^k \rightarrow \forall g^{\mathbb{N} \rightarrow \mathbb{N}} \exists n (|P|^{gn} \rightarrow |P|^n) \\
& \rightsquigarrow \forall g, k \exists n (|P|^k \rightarrow |P|^{gn} \rightarrow |P|^n) \\
& \rightsquigarrow \forall g, k (|P|^k \rightarrow |P|^{gk} \rightarrow |P|^k)
\end{aligned}$$

$$\frac{\begin{array}{ccc} [ |P|^k ] & & [ \neg |P|^{g^0} ] \\ & \vdots & \vdots \\ \exists k |P|^k \vee \forall k \neg |P|^k & |P|^{g^k} \rightarrow |P|^k & |P|^{g^0} \rightarrow |P|^0 \end{array}}{\neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n)}$$

First branch:

$$\begin{aligned}
& \exists k |P|^k \rightarrow \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n) \\
\rightsquigarrow & \exists k |P|^k \rightarrow \forall g^{\mathbb{N} \rightarrow \mathbb{N}} \exists n (|P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g, k \exists n (|P|^k \rightarrow |P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g, k (|P|^k \rightarrow |P|^{g^k} \rightarrow |P|^k)
\end{aligned}$$

Second branch:

$$\begin{aligned}
& \forall k \neg |P|^k \rightarrow \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n) \\
\rightsquigarrow & \forall k \neg |P|^k \rightarrow \forall g \exists n (|P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g \exists k, n (\neg |P|^k \rightarrow |P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g (\neg |P|^{g^0} \rightarrow |P|^{g^0} \rightarrow |P|^0)
\end{aligned}$$

$$\frac{\begin{array}{ccc} [ |P|^{g^0} ] & & [ \neg |P|^{g^0} ] \\ \vdots & & \vdots \\ |P|^{g^0} \vee \neg |P|^{g^0} & |P|^{g(g^0)} \rightarrow |P|^{g^0} & |P|^{g^0} \rightarrow |P|^0 \end{array}}{\neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n)}$$

First branch:

$$\begin{aligned}
& \exists k |P|^k \rightarrow \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n) \\
\rightsquigarrow & \exists k |P|^k \rightarrow \forall g^{\mathbb{N} \rightarrow \mathbb{N}} \exists n (|P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g, k \exists n (|P|^k \rightarrow |P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g, k (|P|^k \rightarrow |P|^{g^k} \rightarrow |P|^k)
\end{aligned}$$

Second branch:

$$\begin{aligned}
& \forall k \neg |P|^k \rightarrow \neg \neg \exists n \forall m (|P|^m \rightarrow |P|^n) \\
\rightsquigarrow & \forall k \neg |P|^k \rightarrow \forall g \exists n (|P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g \exists k, n (\neg |P|^k \rightarrow |P|^{g^n} \rightarrow |P|^n) \\
\rightsquigarrow & \forall g (\neg |P|^{g^0} \rightarrow |P|^{g^0} \rightarrow |P|^0)
\end{aligned}$$

$$\frac{
\begin{array}{ccc}
& [|P|^{g0}] & [\neg|P|^{g0}] \\
& \vdots & \vdots \\
|P|^{g0} \vee \neg|P|^{g0} & |P|^{g(g0)} \rightarrow |P|^{g0} & |P|^{g0} \rightarrow |P|^0
\end{array}
}{
|P|^{g(Ng)} \rightarrow |P|^{Ng}
}$$

Solved by

$$Ng := \begin{cases} 0 & \text{if } \neg|P|^{g0} \\ g0 & \text{if } |P|^{g0} \end{cases}$$

FINITARY DRINKER'S PARADOX II: For an arbitrary function  $g: \mathbb{N} \rightarrow \mathbb{N}$  there exists some  $N: (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \mathbb{N}$  satisfying

$$|P|^{g(Ng)} \rightarrow |P|^{Ng}.$$

This can be defined as

$$Ng := \begin{cases} 0 & \text{if } \neg |P|^{g0} \\ g0 & \text{otherwise.} \end{cases}$$

FINITARY DRINKER'S PARADOX II: For an arbitrary function  $g: \mathbb{N} \rightarrow \mathbb{N}$  there exists some  $N: (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \mathbb{N}$  satisfying

$$|P|^{g(Ng)} \rightarrow |P|^{Ng}.$$

This can be defined as

$$Ng := \begin{cases} 0 & \text{if } \neg |P|^{g0} \\ g0 & \text{otherwise.} \end{cases}$$

The proof theoretic environment is represented by an explicit ‘counterexample function’  $g$ . Any instance of the drinkers paradox in a bigger proof will involve a concrete instantiation  $g_v$  of  $g$ :

$$\begin{aligned} & \exists n \forall m (|P|^m \rightarrow |P|^n) \rightarrow B \\ & \exists U, g \forall n, v (|P|^{gn} \rightarrow |P|^n) \rightarrow |B|_v^{Un} \end{aligned}$$

and hence  $\forall v |B|_v^{U(Ng)v}$  holds.

GENERAL FINITARY DRINKER'S PARADOX: There exists an approximate witness  $\mathcal{N}$  to  $\exists n \forall m (|P|^n \rightarrow |P|^m)$ , that works relative to any environment  $\mathcal{M}$  (representing  $\forall m$ ).

Technique	$\mathcal{N}$	$\mathcal{M}$
$\epsilon$ -calculus	$\begin{cases} \epsilon_n := 0 \\ \text{Check }  P ^{\epsilon_n} \rightarrow  P ^0. \text{ If true, END.} \\ \text{Else } \epsilon_n := \epsilon_m 0 \end{cases}$	$\epsilon_m(\cdot)$
Dialectica	$Ng := \begin{cases} 0 & \text{if } \neg  P ^{g0} \\ g0 & \text{otherwise} \end{cases}$	$g: \mathbb{N} \rightarrow \mathbb{N}$

They both carry out learning, but in completely different frameworks:

$\epsilon$ -calculus  $\rightsquigarrow$  EXPLICITLY via  $\epsilon$ -elimination procedure

Dialectica  $\rightsquigarrow$  IMPLICITLY via contractions in negative-translated proofs



A **learning algorithm**  $\mathcal{L} := (\text{Good}, q, p, \oplus)$  of type  $X, Y, Z$  consists of the following components:

- A decidable predicate **Good** on  $X \times Y$  with the meaning

$\text{Good}(x, y) \Leftrightarrow x$  is good relative to a counterexample object  $y$ ;

- A function  $q: X \rightarrow Z$  which assigns a final value  $q(x) \in Z$  to a good element  $x \in X$ ;
- A function  $p: X \rightarrow Y$  that takes each candidate  $x \in X$  and returns a counterexample object  $p(x) \in Y$ ;
- An operation  $\oplus: X \times Y \rightarrow X$  which updates  $x \in X$  with new information obtained from a counterexample object  $y$ : the update

$$x \mapsto x \oplus y$$

is an instance of *learning*.

The **learning procedure** starting from  $x_0 \in X$  is the sequence defined by

$$\neg \text{Good}(x_i, p(x_i)) \Rightarrow x_{i+1} := x_i \oplus p(x_i).$$

This sequence is finite iff it eventually arrives at some good element  $x_k$  i.e.  $Q(x_k, p(x_k))$ . In this case we say that  $q(x_k)$  is the **limit** of  $\mathcal{L}_{x_0}$  and write

$$\lim_{x_0} \mathcal{L} := q(x_k).$$

REMARK. If  $x \oplus p(x) \prec x$  for some well-founded ordering  $\prec$  then  $\lim_{x_0} \mathcal{L}$  is definable by well-founded recursion on  $\prec$ .

What are learning procedures designed to compute? Suppose that

- $A$  is an initial predicate over  $X$ ,
- $B$  is a target predicate over  $Z$ , and

$$\forall x \begin{cases} \text{Good}(x, p(x)) \rightarrow A(x) \rightarrow B(q(x)) \\ \neg \text{Good}(x, p(x)) \rightarrow A(x) \rightarrow A(x \oplus p(x)) \end{cases}$$

then

$$\forall x_0 (A(x_0) \rightarrow B(\lim_{x_0} \mathcal{L}))$$

whenever the limit exists.

The next two slides are inspired by (Schwichtenberg 2004). Consider the least element principle for  $\Sigma_1$ -formulas relative to some arbitrary well-founded ordering  $\prec$ :

$$\exists x|C|^x \rightarrow \exists y(|C|^y \wedge \forall z \prec y \neg |C|^z).$$

The functional interpretation of this is given by

$$\begin{aligned} &\rightsquigarrow \exists x|C|^x \rightarrow \neg\neg\exists y(|C|^y \wedge \forall z \prec y \neg |C|^z) \\ &\rightsquigarrow \exists x|C|^x \rightarrow \forall p\exists y(|C|^y \wedge (p(y) \prec y \rightarrow \neg |C|^{p(y)})) \\ &\rightsquigarrow \forall x, p\exists y(|C|^x \rightarrow |C|^y \wedge (p(y) \prec y \rightarrow \neg |C|^{p(y)})) \quad (*) \end{aligned}$$

THEOREM. Define  $\mathcal{L}_p := (\text{Good}_C, \text{id}, p, \pi_1)$  for

$$\text{Good}_C(y, z) :\equiv z \prec y \rightarrow \neg |C|^z.$$

Then (\*) is realized by

$$\lambda x, p . \lim_x \mathcal{L}_p$$

For any  $a \geq b > 0$  there exist some  $n, m$  such that  $am + bn \mid a, b$

Proof. Define the measure  $\mu: \mathbb{N}^2 \rightarrow \mathbb{N}$  by  $\mu \underline{x} := \underline{x} \cdot (a, b)$  and the ordering  $\succ$  on  $(\mathbb{N}^2)^*$  by

$$\square \succ [\underline{x}] \text{ and } s * \underline{x} \succ s * \underline{x} * \underline{y} \text{ iff } \mu \underline{x} > \mu \underline{y}.$$

Define  $|C|^s := \forall i < |s| (\mu \tilde{s}_{i+2} > 0 \wedge \mu \tilde{s}_{i+2} = \text{rem}(\mu \tilde{s}_i, \mu \tilde{s}_{i+1}))$ , where  $\tilde{s} := [e_0, e_1] * s$ .

Then  $|C|^\square$  trivially holds, and thus there is some minimal  $s$  satisfying  $|C|^s$ .

Suppose that  $\mu \tilde{s}_l = q \mu \tilde{s}_{l+1} + r$  and define  $\underline{r} = \tilde{s}_l - q \tilde{s}_{l+1}$ . Then  $s \succ s * \underline{r}$ , but  $\neg |C|^{s * \underline{r}}$  implies that  $\mu \underline{r} = 0$  and thus  $\mu \tilde{s}_{l+1} \mid \mu \tilde{s}_l$ .

By induction we must have  $\mu \tilde{s}_{l+1} \mid \mu \tilde{s}_i$  for all  $i \leq l$  and thus  $\mu \tilde{s}_{l+1} \mid \mu e_0, \mu e_1$  i.e.  $\mu \tilde{s}_{l+1} \mid a, b$ .

What is the computational content of this proof? The main non-constructive component is a single instance of the least element principle, so should essentially be a simple learning procedure modulo some details.

The formula  $\forall a, b \exists m, n (am + bn \mid a, b)$  is realized by the program

$$\lambda a, b . \lim_{\square} \mathcal{L}_{a,b}$$

where  $\mathcal{L}_{a,b} := (\text{Good}, q, p_{a,b}, *)$  for

- $p_{a,b}(s) = \tilde{s}_l - \text{quot}(\mu_{a,b}\tilde{s}_l, \mu_{a,b}\tilde{s}_{l+1})$  and  $q = \tilde{s}_{l+1}$ , and
- $\text{Good}_{a,b}(s, t) := t \prec s \rightarrow \neg |C|^t$

CLAIM:  $\mathcal{L}_{a,b}$  is just the Euclidean algorithm!

A core part of most learning-based computational interpretations of classical logic involves understanding how to deal with the so called  $\Sigma_1\text{-LEM}^-$ :

$$\forall n(\exists x|P_n|^x \vee \forall y\neg|P_n|^y).$$

Much work has been done recently on extracting programs from proofs relative to non-computable Skolem functions  $f_P$  for  $\Sigma_1\text{-LEM}^-$  satisfying

$$\exists f \forall n, y(|P_n|^{f^n} \vee \neg|P_n|^y).$$

General idea behind  $\epsilon$ -calculus and more modern approaches (e.g. ‘interactive learning-based realizability’ of Aschieri and Berardi’):

1. If  $\text{HA} + \Sigma_1\text{-LEM}^- \vdash A$  then there is a realizing term  $t$  for  $A$  that is primitive recursive in Skolem functions  $f_P$ .
2. When  $A$  is e.g.  $\Pi_2$  the term  $t$  only requires *approximations* to  $f_P$ . Therefore an approximation  $t_0 \sqsubset t_1 \sqsubset \dots \sqsubset t$  to  $t$  can be built by ‘learning’ a finite amount of information about  $f_P$ .

An equivalent problem involves giving a computational interpretation to the axiom of countable choice for  $\Pi_1^0$ -formulas:

$$\forall n \exists x \forall y |A_n|_y^x \rightarrow \exists f \forall n, y |A_n|_y^{f(n)}$$

This has a functional interpretation given by

$$\begin{aligned} &\rightsquigarrow \forall n \neg \neg \exists x \forall y |A_n|_y^x \rightarrow \neg \neg \exists f \forall n, y |A_n|_y^{f(n)} \\ &\rightsquigarrow \forall n, p \exists x |A_n|_{p(x)}^x \rightarrow \forall \varphi, \phi \exists f |A_{\varphi f}|_{\phi f}^{f(\varphi f)} \\ &\rightsquigarrow \forall L, \varphi, \phi \exists f (\forall n, p |A_n|_{p(Lnp)}^{Lnp} \rightarrow |A_{\varphi f}|_{\phi f}^{f(\varphi f)}) \end{aligned}$$

Now suppose that the premise is realized by a countable sequence  $(\mathcal{L}[n])_{n \in \mathbb{N}}$  of learning procedures i.e.

$$Lnp := \lim_{x_n} \mathcal{L}[n]_p$$

for  $\mathcal{L}[n]_p := (\text{Good}_n, q_n, p \circ q_n, \oplus_n)$  and  $(x_n)_{n \in \mathbb{N}}$  some sequence of initial objects.



We define a ‘product’  $\mathcal{L}[\infty]$  of the  $(\mathcal{L}[n])_{n \in \mathbb{N}}$ , namely

$$\mathcal{L}[\infty]_{\varphi, \phi} := (\mathbf{Good}_{\infty}, q_{\infty}, [\varphi, \phi] \circ q_{\infty}, \oplus_{\infty})$$

has type  $X^{\mathbb{N}}, \mathbb{N} \times X$  and

- $f \oplus_{\infty} (n, x) := f[n \mapsto f(n) \oplus_n x]$ ;
- $\mathbf{Good}_{\infty}(f, [n, x]) := T_n(f(n), x)$ ; and
- $q_{\infty}(f) := \lambda n. q_n(f(n))$ .

**THEOREM.** The formula

$$\forall \varphi, \phi \exists f (\forall n, p |A_n|_p^{\lim_{x_n} \mathcal{L}[n]_p} \rightarrow |A_{\varphi f}|_{\phi f}^{f(\varphi f)})$$

is realised by

$$\lambda \varphi, \phi . \lim_{f_0} \mathcal{L}[\infty]_{\varphi, \phi}$$

for  $f_0 := \lambda n. x_n$ .

Functional interpretation of  $\neg\neg(\exists x|P_n|^x \vee \forall y\neg|P_n|^y)$  given by

$$\forall p \exists b, x(|P_n|^x \vee_b \neg|P_n|^{p(x)})$$

and is realised by  $\lambda p. \lim_{(\perp, 0)} \mathcal{L}[n]_p$  for  $\mathcal{L}[n]_p := (\mathbf{Good}_n, \text{id}, p, \pi_1)$  with

$$\mathbf{Good}_n((b, x), (b', y)) := b = \perp \rightarrow \neg|P_n|^y.$$

Therefore the functional interpretation of  $\neg\neg\exists f \forall n(|P_n|^{f(n)} \vee \forall y\neg|P_n|^y)$

$$\forall \varphi, \phi \exists f(|P_{\varphi f}|^{f_1(\varphi f)} \vee_{f_0(\varphi f)} \neg|P_{\varphi f}|^{\phi f})$$

is realised by  $\lambda \varphi, \phi. \lim \mathcal{L}[\infty]_{\varphi, \phi}$  for  $\mathcal{L}[\infty]_{\varphi, \phi} := (\mathbf{Good}_\infty, \text{id}, [\varphi, \phi], \oplus_\infty)$  where

- $f \oplus_\infty (n, x) := f[n \mapsto x]$  is a standard function update operation and
- $\mathbf{Good}_\infty(f, [n, (b, x)]) \Leftrightarrow n \notin \text{dom}(f) \rightarrow \neg|P_n|^x$

## DIRECTIONS FOR RESEARCH:

- 1 How does this compare to traditional interpretation of choice principles via recursors? Hopefully (P. 2016)!
- 2 Program extraction from proofs in mathematics: many theorems in subsystems of analysis use only  $\Pi_1^0$ -comprehension e.g. Ramsey's theorem (Kreuzer, Kohlenbach 2009). Can we give a learning based interpretation of these? What about more complex things like Kruskal's theorem?
- 3 Can we obtain complexity results analogous to those for epsilon calculus and interactive learning realizability? Potentially show that limits are computable using some kind of bar recursion of low type and adapt (Schwichtenberg 1979)...

THANK YOU!