

# On fixed points, diagonalization, and self-reference

Bernd Buldt

Department of Philosophy  
Indiana U - Purdue U Fort Wayne (IPFW)  
Fort Wayne, IN, USA  
e-mail: buldtb@ipfw.edu

CL 16 – Hamburg – September 12, 2016

# Section I: G1 & Fixed Points

# G1 Proof, using the Gödel fixed point

## Assumptions

(ADQ)  $\vdash_{\mathcal{F}} \varphi \Leftrightarrow \vdash_{\mathcal{F}} \text{Pr}_{\mathcal{F}}(\ulcorner \varphi \urcorner)$ , for all  $\varphi \in \mathcal{L}_{\mathcal{F}}$

(FPE)  $\vdash_{\mathcal{F}} \gamma \Leftrightarrow \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner)$ , for at least one  $\gamma \in \mathcal{L}_{\mathcal{F}}$

## Proof

$$\vdash_{\mathcal{F}} \gamma \stackrel{\text{ADQ}}{\Rightarrow} \vdash_{\mathcal{F}} \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner) \stackrel{\text{FPE}}{\Rightarrow} \vdash_{\mathcal{F}} \neg \gamma \Rightarrow \not\vdash_{\mathcal{F}} \gamma \stackrel{\text{con}_{\mathcal{F}}}{\Rightarrow} \not\vdash_{\mathcal{F}} \gamma$$

$$\vdash_{\mathcal{F}} \neg \gamma \stackrel{\text{FPE}}{\Rightarrow} \vdash_{\mathcal{F}} \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner) \stackrel{\text{ADQ}}{\Rightarrow} \vdash_{\mathcal{F}} \gamma \Rightarrow \not\vdash_{\mathcal{F}} \neg \gamma \stackrel{\text{con}_{\mathcal{F}}}{\Rightarrow} \not\vdash_{\mathcal{F}} \neg \gamma$$

## Fixed point derivation, Step 1: Substitution

- ▶ Fix a certain individual variable of your choice; say 'u.'
- ▶ Define a function *sub* that mirrors the substitution of the replacee variable 'u' for a replacer term 't,'

$$\varphi[u]_u^t \equiv \varphi(t),$$

but in the realm of Gödel numbers. In short:

$$\text{sub}(x, y) := \begin{cases} \text{gn}(\varphi[u]_u^{\bar{t}}) & \text{if } x = \text{gn}(\varphi(u)) \text{ and } y = \text{gn}(\bar{t}) \\ x & \text{otherwise.} \end{cases}$$

- ▶ Note that  $\text{sub}(x, y)$  is primitive recursive and therefore represented by an expression  $\varphi_s(x, y)$  in  $\mathcal{F}$ .

## Fixed point derivation, Step 2: Definitions

- ▶ Define  $\varphi(u) := \forall x [\neg \text{Proof}_F(x, \text{sub}(u, u))]$ .
- ▶ Define  $p := gn(\varphi(u))$ .
- ▶ Substitute  $p$  for  $u$  in  $\varphi(u)$ , viz.,

$$\gamma := \varphi(\bar{p}) \equiv \forall x [\neg \text{Proof}_F(x, \text{sub}(\bar{p}, \bar{p}))].$$

- ▶ Calculate
 

$sub(p, p)$	$=$	$sub(gn(\varphi(u)), p)$	; def. $p$
	$=$	$gn(\varphi[u]_{\bar{p}}^u)$	; def. $sub$
	$=$	$gn(\varphi(\bar{p}))$	; substitution
	$=$	$gn(\gamma)$	; def. $\gamma$

## Fixed point derivation, Step 3: Derivation

- ▶ Recall Step 2:  $sub(p, p) = gn(\gamma)$ .

- ▶ Reason inside  $\mathcal{F}$ .

$$\begin{array}{ll} \vdash_{\mathcal{F}} \neg Pr_{\mathcal{F}}(x) \leftrightarrow \neg Pr_{\mathcal{F}}(x) & ; \text{ logic} \\ \vdash_{\mathcal{F}} \neg Pr_{\mathcal{F}}(sub(\bar{p}, \bar{p})) \leftrightarrow \neg Pr_{\mathcal{F}}(\ulcorner \gamma \urcorner) & ; \text{ Step 2} \\ \vdash_{\mathcal{F}} \forall x [\neg Proof_{\mathcal{F}}(x, sub(\bar{p}, \bar{p}))] \leftrightarrow \neg Pr_{\mathcal{F}}(\ulcorner \gamma \urcorner) & ; \text{ def. } Pr_{\mathcal{F}} \\ \vdash_{\mathcal{F}} \varphi(\bar{p}) \leftrightarrow \neg Pr_{\mathcal{F}}(\ulcorner \gamma \urcorner) & ; \text{ def. } \varphi(\bar{p}) \\ \vdash_{\mathcal{F}} \gamma \leftrightarrow \neg Pr_{\mathcal{F}}(\ulcorner \gamma \urcorner) & ; \text{ def. } \gamma \end{array}$$

- ▶ Warning. We assumed  $\vdash_{\mathcal{F}} sub(\bar{p}, \bar{p}) = \ulcorner \gamma \urcorner$ , which requires induction.

## Theorem (Fixed Point Theorem, Diagonalization Lemma)

*Assume  $\mathcal{F}$  to allow for representation. For each expression  $\varphi$  with at least one variable free, there is a  $\psi$  such that,*

$$\vdash_{\mathcal{F}} \psi \leftrightarrow \varphi_{\psi}$$

*where  $\varphi_{\psi}$  can be either of the four forms:*

$$\varphi(\ulcorner \psi \urcorner), \varphi(\ulcorner \neg \psi \urcorner), \neg \varphi(\ulcorner \psi \urcorner), \neg \varphi(\ulcorner \neg \psi \urcorner),$$

*viz., instances of what we call a Henkin, Jeroslov, Gödel, or Rogers fixed point resp.*

## Proof.

Same as above (with minor modifications). □

## Black self-referential magic?

- ▶ Two questions about fixed points such as

$$\vdash_{\mathcal{F}} \gamma \leftrightarrow \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner).$$

1. How much “black magic” is required for their derivation?  
... will be answered in Section II.
2. How much “self-reference” do they involve?  
... will be answered in Section III.



# Section II: Diagonalization

# Black magic?

## 1<sup>st</sup> Question

How much “black magic” is required for the derivation of fixed points such as

$$\vdash_{\mathcal{F}} \gamma \leftrightarrow \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner)?$$

## Answer

None.

## Diagonalization

- ▶ Let  $\mathcal{A} = \{a_{ij}\}_{i,j \in \omega}$  be a (countable) two-dimensional array:

$$\begin{array}{cccccc} R_0 : & a_{00} & a_{01} & \dots & a_{0n} & \dots \\ R_1 : & a_{10} & a_{11} & \dots & a_{1n} & \dots \\ & \vdots & \vdots & \ddots & \vdots & \\ R_n : & a_{n0} & a_{n1} & \dots & a_{nn} & \dots \\ & \vdots & \vdots & & \vdots & \ddots \end{array}$$

- ▶ Let  $f$  be a sequence transforming function,

$$f(R_n) = \{f(a_{ni})\}_{i \in \omega}.$$

- ▶ Apply  $f$  to the diagonal sequence  $D$ :

$$D' = f(D) := \langle f(a_{00}), f(a_{11}), f(a_{22}), \dots, f(a_{nn}), \dots \rangle.$$

## Diagonalization: (Non-)Closure

- ▶ One of two things can happen to the anti-diagonal  $D' = f(D)$ :
  1.  $D'$  is identical to one of the rows, *viz.*,  $f(D) = R_i \in \mathcal{A}$ , for some  $i$ .
  2.  $D'$  is not identical to any of the rows, *viz.*,  $f(D) \neq R_i \in \mathcal{A}$ , for all  $i$ .
- ▶ If Case 1 applies, we call the set  $A$  closed under  $f$ , and  $f$  will have fixed points.
- ▶ If Case 2 applies,  $A$  is not closed under  $f$ , and we have Cantor's diagonal argument showing that a certain sequence is not in  $\mathcal{A}$  (to “diagonalize out”).

## Diagonalization: Case 1 – Closure

- ▶  $D'$  is identical to one of the rows, *viz.*,  $f(D) = R_i \in \mathcal{A}$ , for some  $i$ .

- ▶ The identity  $D' = f(D) = R_i$  is element-wise identity:

$$\begin{array}{ccccccccc}
 D' & = & \langle f(a_{00}), & f(a_{11}), & \dots, & f(a_{ii}), & \dots, & f(a_{nn}), & \dots \rangle \\
 & & \parallel & \parallel & & \parallel & & \parallel & \\
 R_i & = & \langle a_{i0}, & a_{i1}, & \dots, & a_{ii}, & \dots, & a_{in}, & \dots \rangle
 \end{array}$$

- ▶ Closure under  $f$  (failure to “diagonalize out” ) implies fixed points  $f(a_{ij}) = a_{ij}$ .

## Diagonalization: Case 1 – Closure

$$\begin{array}{cccccc}
 R_0 : & a_{00} & a_{01} & \dots & a_{0n} & \dots \\
 R_1 : & a_{10} & a_{11} & \dots & a_{1n} & \dots \\
 & \vdots & \vdots & \ddots & \vdots & \\
 R_n : & a_{n0} & a_{n1} & \dots & a_{nn} & \dots \\
 & \vdots & \vdots & & \vdots & \ddots
 \end{array}
 \Rightarrow
 \begin{array}{cccccc}
 R_0 : & fa_{00} & a_{01} & \dots & a_{0n} & \dots \\
 R_1 : & a_{10} & fa_{11} & \dots & a_{1n} & \dots \\
 & \vdots & \vdots & \ddots & \vdots & \\
 R_n : & a_{n0} & a_{n1} & \dots & fa_{nn} & \dots \\
 & \vdots & \vdots & & \vdots & \ddots
 \end{array}$$

$$\Rightarrow
 \begin{array}{cccccc}
 R_0 : & a_{00} & a_{01} & \dots & a_{0i} & \dots & a_{0n} & \dots \\
 R_1 : & a_{10} & a_{11} & \dots & a_{1i} & \dots & a_{1n} & \dots \\
 & \vdots & \vdots & \ddots & \vdots & & \vdots & \\
 \Rightarrow & f(D) = R_i : & \frac{fa_{00}}{a_{i0}} & \frac{fa_{11}}{a_{i1}} & \dots & \frac{fa_{ij}}{a_{ij}} & \dots & \frac{fa_{nn}}{a_{in}} & \dots \\
 & & \vdots & \vdots & & \vdots & \ddots & \vdots & \\
 R_n : & a_{n0} & a_{n1} & \dots & a_{ni} & \dots & a_{nn} & \dots
 \end{array}$$

## Diagonalization: Closure & Gödel fixed point

- ▶ Can we understand  $\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner)$  to be an instance of  $f(a_{ii}) = a_{ii}$  for some  $f$  and some array  $\mathcal{A} = \{a_{ij}\}_{i,j \in \omega}$ ?
- ▶ Yes.

## Diagonalization: Closure & Gödel fixed points

- ▶ Step 1: Choose all first-order expressions with the free variable 'u:'

$$A = \{\varphi_0(u), \varphi_1(u), \varphi_2(u), \dots\}.$$

- ▶ Step 2: Form the set of all of their Gödel numbers:

$$B = \{\ulcorner \varphi_0(u) \urcorner, \ulcorner \varphi_1(u) \urcorner, \ulcorner \varphi_2(u) \urcorner, \dots\}.$$

- ▶ Step 3: Systematically plug all members of  $B$  into the free variable slots of all members of  $A$ ; call this set  $C$ . We write ' $\varphi_{ab}$ ' instead of ' $\varphi_a(\ulcorner \varphi_b \urcorner)$ '.



## Diagonalization: Gödel fixed points – 1<sup>st</sup> diagonalization

- Lay out the elements of  $C$  in such a way that  $A$  determines the rows and  $B$  the columns which gives us::

	$\ulcorner \varphi_0 \urcorner$	$\ulcorner \varphi_1 \urcorner$		$\ulcorner \varphi_n \urcorner$	
$\varphi_0$	$\varphi_{00}$	$\varphi_{01}$	$\dots$	$\varphi_{0n}$	$\dots$
$\varphi_1$	$\varphi_{10}$	$\varphi_{11}$	$\dots$	$\varphi_{1n}$	$\dots$
	$\vdots$	$\vdots$	$\ddots$	$\vdots$	
$\varphi_n$	$\varphi_{n0}$	$\varphi_{n1}$	$\dots$	$\varphi_{nn}$	$\dots$
	$\vdots$	$\vdots$		$\vdots$	$\ddots$

- Note that the diagonal sequence  $\{\varphi_{xx}\}_{x \in \omega}$  corresponds to the substitution function  $sub(x, x)$  we used above.

## Diagonalization: Gödel fixed points – 2<sup>nd</sup> diagonalization

1. Observe that the provability predicate  $\neg\text{Pr}_F(u)$  is itself part of the first set we started out with:  $A = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$ ; i. e.,  $\exists i$  s. t.:  $\varphi_i \equiv \neg\text{Pr}_F(u)$ .
2. Apply the transformation  $f : \varphi_{ab} \mapsto \neg\text{Pr}_F(\varphi_{ab})$ .
3. Because of (1),  $f$  maps  $C$  onto  $C$ ,  $C$  will be closed under  $f$ , and each image  $\neg\text{Pr}_F(\varphi_{ab})$  must be a  $\varphi_{in}$ , for some  $n$ .
4. Hence,  $f(D)$  has a fixed point  $\varphi_{ii}$ , which corresponds to the expression  $\gamma \equiv \varphi(\bar{p})$  we used above.

## Diagonalization: Gödel fixed points without “black magic”

- ▶ Derivable fixed points in systems of arithmetic  $\mathcal{F}_{Ar}$ , e. g.,

$$\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner),$$

are a result of the fact that set of expressions, such as  $A$ , are closed under certain transformations  $f$ .

- ▶  $sub(x, x)$  corresponds to  $\{\varphi_{xx}\}_{x \in \omega}$ .
- ▶  $\gamma \equiv \varphi(\bar{p})$  corresponds to  $\varphi_{ii}$ .
- ▶ Outcomes can be modelled in  $\mathcal{F}_{Ar}$ .
- ▶ The procedure (“double diagonalization”) is entirely syntactic is completely mundane, no magic anywhere.

# Section III: Self-Reference

# Black magic?

## 2<sup>nd</sup> Question

How much “self-reference” is required for the derivation of fixed points such as:

$$\vdash_{\mathcal{F}} \gamma \leftrightarrow \neg \text{Pr}_{\mathcal{F}}(\ulcorner \gamma \urcorner)?$$

## Answer

None.

## Self-Reference: Rendered moot by diagonalization

- ▶ Previous section: Fixed points such as:

$$\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner),$$

result from certain closure properties.

- ▶ The crucial steps,
  - ▶  $\text{sub}(x, x)$  or  $\{\varphi_{xx}\}_{x \in \omega}$ .
  - ▶  $\gamma \equiv \varphi(\bar{p})$  or  $\varphi_{ii}$ .

are entirely syntactic operations, which neither employ nor presuppose any concept of self-reference.

## Self-Reference: Digging deeper

- ▶ Does  $\psi \leftrightarrow \varphi(\psi)$  mean that  $\psi$  says it has property  $\varphi$ ?
  - ▶ Does  $\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner)$  mean that  $\gamma$  expresses some property it itself has, namely, the property “ $\neg \text{Pr}_F(u)$ ” (unprovability)?
  - ▶ If so, does it mean that  $\gamma$  states its own unprovability?
- ▶ Preliminaries: What self-reference cannot be.
  - ▶ Self-reference cannot mean  $\gamma$  is somehow a proper part of itself; this would violate the mereological definition of proper parthood,  $PP_{xy} := P_{xy} \wedge x \neq y$ .
  - ▶ Self-reference hence presupposes a more abstract semantical relation than self-inclusion is.

## Self-Reference: 'Propertual' self-reference

- ▶ Expression  $\varphi(u)$  defines, in some structure  $\mathfrak{A}$ , property  $P$  if:

1. Definition:  $\{x : P(x)\}$  iff  $\{x : \mathfrak{A} \models \varphi(\#x)\}$ .

Then  $\varphi(u)$  has property  $P$  itself if:

2. Self-Reference:  $\mathfrak{A} \models \varphi(\#\varphi(u))$ .

- ▶ Application to  $\neg\text{Pr}_F(u)$

- ▶  $\mathfrak{A} \models \neg\text{Pr}_F(\ulcorner\neg\text{Pr}_F(u)\urcorner)$ , because  $\not\models_{\mathcal{F}} \neg\text{Pr}_F(u)$

- ▶ Given suitable circumstances, 'propertual' self-reference may occur.

- ▶ Mute point: no mention of  $\gamma \leftrightarrow \neg\text{Pr}_F(\ulcorner\gamma\urcorner)$ .



## Self-Reference: Propertual self-reference

- ▶ Problem. What conditions would elevate  $\psi$  in  $\psi \leftrightarrow \varphi_\psi$  from being merely truth-functionally equivalent to actually being self-referential the same way  $\varphi_\psi$  is?
- ▶ All known attempts to identify such conditions can be considered to have failed, mostly because we do not yet have a good theory of self-reference.  
(see Halbach and Visser 2015)

## Self-Reference: Improper self-reference

Direct objectual self-reference:  $\varphi(\#\varphi)$ ; eg, viz.,  $\varphi \wedge |\varphi|$ , or  $\varphi(\ulcorner \varphi \urcorner)$ .

- ▶ Does  $\gamma$  in  $\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner)$  contain its own name?
- ▶ Recall that  $\gamma$  is shorthand for  $\forall x[\neg \text{Proof}_F(x, \text{sub}(\bar{p}, \bar{p}))]$ , with  $p = gn(\neg \text{Pr}_F(\text{sub}(u, u)))$ .
- ▶ Thus, no.
- ▶ However, since  $\text{sub}(\bar{p}, \bar{p}) = gn(\gamma)$ , we know that  $\gamma$  would be self-referential if criteria would be more lax.

## Self-Reference: Improper self-reference

Indirect objectual self-reference:  $\varphi(\#\#\varphi)$ ; eg,  $\varphi(t)$ , with  $t = \#\#\varphi(t)$

- ▶ Does  $\gamma$  in  $\gamma \leftrightarrow \neg \text{Pr}_F(\ulcorner \gamma \urcorner)$  contain its own indirect name?
- ▶ Since  $\text{sub}(\bar{p}, \bar{p}) = \text{gn}(\gamma)$ , the expression  $\gamma$ , which is  $\forall x[\neg \text{Proof}_F(x, \text{sub}(\bar{p}, \bar{p}))]$ , contains an indirect name of itself.
- ▶ Some (eg, Heck 2007) are perfectly happy to embrace the last point and call the Gödel sentence  $\gamma$  self-referential in the above sense and have it say “I’m not provable.”

## Self-Reference: Improper self-reference

- ▶  $\gamma$  does not say “I” but refers to itself indirectly via a functional expression
- ▶  $\gamma$  is true *iff*  $\gamma$  is not formally provable. By itself, this is a raw datum about  $\gamma$ 's model theoretic evaluation and the resulting truth value. As such, it is just another equivalence that implies nothing about meaning or self-reference.
- ▶ Semantic stance like intentional stance; useful but not justified
- ▶ We practice semantic hunches, but gut feelings are a poor substitute for an actual theory.

## Self-Reference: Summary

- ▶ Diagonalization produces fixed points.
- ▶ Fixed points do not establish self-reference.
- ▶ Self-reference we find is not proper internal self-reference, but our external attribution.

Thank You!