

# The anti-mechanism argument based on Gödel's incompleteness theorems, indescribability of the concept of natural number and deviant encodings

Paula Quinon

Warsaw University of Technology

Department of Administration and Social Sciences

International Center for Formal Ontology

---

## Abstract

This paper reassesses the criticism of *the Lucas-Penrose anti-mechanism argument based on Gödel's incompleteness theorems* formulated by Krajewski [17]. Krajewski claims that in order for *the argument* to work, *an extra-formal assumption* should be added: “the human mind is consistent”. The assumption is “extra-formal”, because it cannot be formalised and must be accepted on the basis of a non-formal insight. Thus – claims Krajewski – the Lucas-Penrose anti-mechanism argument, which relies on the formalisation of the whole process, fails to establish that human mind is not mechanistic.

An additional light on Krajewski's rejection of the anti-mechanism argument, is shed by *a corollary to this argument*: the human mind allegedly outperforms machines, because although there is no purely formal, exhaustive definition of natural numbers, human mathematicians can successfully work with natural numbers. Also in this case, as Krajewski claims, the anti-mechanism argument requires an additional assumption, such as “PA1 is complete” or “there exists a set of all natural numbers”, which must be accepted without a full formal explanation of its meaning. Again, the anti-mechanism argument, which requires full formalisation, fails.

I agree with Krajewski that additional extra-formal assumptions are necessary to make the anti-mechanism argument, and its corollary, workable. However, I do not agree with Krajewski's statement that assumptions cannot be formalised. Instead, I propose an alternative problem with “extra-formal” assumptions that call into question the anti-mechanism argument, namely *circularity*. A machine cannot prove its own consistency; a human

mind, to show that it can outperform a machine in this respect, cannot simply bring in the assumption “I am consistent”. Such a reasoning is circular. Starting with an analysis of circularity, I propose a way of thinking about the interplay between informal and formal resources in mathematical reasoning.

*Keywords:* the Lucas-Penrose anti-mechanism argument, Gödel’s incompleteness theorems, the concept of natural number, the concept of computation, Carnapian explications, conceptual engineering, conceptual fixed points, conceptual vicious circles, deviant encodings, structuralism, computational structuralism, the Church-Turing thesis

---

1	<b>Contents</b>	
2	<b>1 The Lucas-Penrose argument and its criticism</b>	<b>5</b>
3	<b>2 The concept of natural number and the concept of compu-</b>	
4	<b>tation</b>	<b>8</b>
5	<b>3 Nested vicious circles</b>	<b>11</b>
6	<b>4 Conceptual engineering and conceptual fixed points</b>	<b>15</b>
7	<b>5 The Church-Turing Thesis as a Carnapian explication</b>	<b>17</b>
8	<b>6 Theory of mind and computations</b>	<b>20</b>
9	<b>7 The Lucas-Penrose argument and extra-formal concepts</b>	<b>23</b>
10	<b>Introduction</b>	

11 The Lucas-Penrose anti-mechanism argument against computability of  
12 the human mind in a nutshell states what follows. According to Gödel’s  
13 incompleteness theorems, there does not exist a (sufficiently rich) consistent  
14 theory that can prove its own consistency. However, mathematical practice  
15 shows that Gödel-type results are commonly proven by human mathemati-  
16 cians. In consequence, says the argument, human mathematicians are not  
17 describable as formal proof systems, nor are reducible to performing algo-  
18 rithms.

19 In [17], Krajewski criticizes the Lucas-Penrose argument by claiming that  
20 Gödel’s incompleteness theorems *stand alone* (as it is in the Lucas-Penrose  
21 case) are not sufficient for formulating the claim that the human mind is  
22 non-computational. The anti-mechanism argument based on Gödel’s incom-  
23 pleteness theorems needs to be enriched by an extra-formal assumption. For  
24 instance, an assumption that the theory constituting the human mind is  
25 consistent.

26 In order to provide an additional context to his investigations, Krajewski  
27 [17], highlights the analogy between the claim that Gödel’s incompleteness  
28 theorems imply the non-computational nature of the human mind, and the  
29 claim that “we [humans] cannot give a definition of the natural numbers as  
30 we understand them” [17, page 31]. The analogy goes as follows: in order to  
31 make a successful anti-mechanism argument based on Gödel’s incompleteness  
32 theorems, one needs to assume – in addition to the formal counterpart – that  
33 the theory constituting the human mind is consistent. The fact that Gödel’s  
34 argument can be iterated for increasingly rich theories is not sufficient for  
35 formulation of the anti-mechanism argument. The possibility to iterate in-  
36 creasingly rich theories, which all have a Gödel’s sentence, and none of which  
37 proves its own consistency, is a formal process and as such can be executed  
38 with purely formal means. Thus, it does not say anything about computabil-  
39 ity or non-computability of the human mind. In order to be able to formulate  
40 the anti-mechanism argument, one needs to assume – for instance – that the  
41 human mind is consistent. Analogously, each definition of natural number  
42 ends up in a vicious circle of definitions, or – as Krajewski says – “our axioms  
43 [both the first-order (*PA1*) and the second-order Peano Arithmetic (*PA2*)]  
44 define numbers only when taken together with some background knowledge  
45 or apparatus that makes possible our intuitive grasp of numbers [such as  
46 the intuition that the first-order Peano’s Arithmetic is complete or the intu-  
47 ition that there exists the set of all natural numbers being referred to in the  
48 background of the second-order Peano’s Arithmetic]” [17, page 31]. In both  
49 cases, an immediate, but incorrect according to Krajewski, conclusion could  
50 be that “no computer can be taught our concept of a number” and that in  
51 consequence “we [humans] are better than any machine” [17, page 31].

52 In this paper, I observe that this analogy can be pushed further to a cir-  
53 cular reasoning. In both cases, making an extra-formal assumption leads to  
54 a vicious circle because one assumes consistency of one’s mind while proving  
55 that the human mind outperforms machines, or one assumes that the con-  
56 cept of set of natural numbers can be intuitively apprehended while defining

57 natural numbers. There exist studies showing that the method of concep-  
58 tual analysis is particularly sensitive to falling into circular reasoning. The  
59 circularity related to the concept of natural number has been investigated  
60 in discussions about *computational structuralism* (Halbach & Horsten 2005  
61 [14], Quinon & Zdanowski 2007 [36]). Computational structuralism is a po-  
62 sition, according to which the concept of natural number and the concept of  
63 computation are closely related. More precisely, according to this position,  
64 an adequate account of natural numbers treats them as objects that can be  
65 used for computations. After a brief overview of the anti-mechanism argu-  
66 ment and its criticism in **Section 1**, in **Section 2** I explain inter-relation  
67 and inter-definability between the concept of natural number and the con-  
68 cept of computation. In **Section 3**, I describe how the two concepts fall into  
69 a vicious circle of definition individually, and also while used in definition of  
70 one another.

71 Rescorla (2007 [41]) identifies problems with conceptual analysis related  
72 to the concept of computation, Quinon (2018 [38]) suggests that there is  
73 no fully satisfactory way out from vicious circles in definitions within con-  
74 ceptual analysis. Approaching the concept of computation and the concept  
75 of natural number from another methodological perspective, seems to be  
76 more fruitful. For instance, an interesting insight can be gained thanks to  
77 *conceptual engineering*. Both concepts have a form of what in the area of  
78 conceptual engineering is called “conceptual fixed point”. A conceptual fixed  
79 point is an idea issued from conceptual engineering of moral concepts, where  
80 it is claimed that some basic moral concepts should not be engineered, but  
81 should always be understood in the most objective way (Eklund 2015 [10]).  
82 **Section 4** is devoted to presentation of the method of conceptual engineer-  
83 ing and adequacy of conceptual fixed points for the concept of computation  
84 and the concept of natural number. As suggested by the phenomenon of  
85 conceptual fixed points, the only way out from vicious circles consists in an  
86 arbitrary decision what is the intended meaning of the given concept.

87 In **Section 5** I extend my methodological investigations into yet another  
88 method, and I discuss advantages of thinking about formalisation of the  
89 concept of computation in terms of Carnapian explications. It has been  
90 argued, for instance in (Quinon 2019 [39]), that a move from an intuitive  
91 concept of computation, used in everyday life, to a scientific or formal concept  
92 as stated by the Church-Turing thesis, follows the schema of a Carnapian  
93 explication. In **Section 6**, I extend the context of Carnapian explications  
94 of the temporary aspect. I realise that both, the concept of natural number

95 and the concept of computation, have been evolving in such a way, that their  
96 core meanings were shifting. I propose a hypothesis that at least a part  
97 of the confusion regarding the specificity of the conceptual structure of the  
98 concept of computation contributes to the confusion regarding the nature of  
99 human reasoning and the human mind. In consequence, I claim that – at  
100 least partially – the “feeling” that there exist non-computational processes  
101 comes from the complexity of the conceptual structure of the concept of  
102 computation.

103 In the final **Section 7**, I wrap up with the ways in which my observations  
104 regarding the concept of computation and the concept of natural number,  
105 could be used for understanding the reasons for which the anti-mechanism  
106 argument fails. I suggest a different reason from the one proposed by Krajewski,  
107 for which the extra-formal assumption prevents the anti-mechanism argu-  
108 ment from success. Firstly, I claim that thanks to the method of Carnapian  
109 explications, it is highly possible to go from intuitive pre-scientific concept to  
110 a formal concept. Secondly, I observe that the extra-formal assumption after  
111 an arbitrary formalisation, leads to the vicious circle in reasoning. Therein  
112 lies the problem.

## 113 1. The Lucas-Penrose argument and its criticism

114 In this section, I present a brief overview of various versions of the anti-  
115 mechanism argument based on Gödel’s incompleteness theorems, and the  
116 ways in which those arguments were criticised. In particular, I explicate  
117 Krajewski’s way of refuting the argument. In my overview, I prioritise the  
118 authors to which Krajewski refers in his paper.

119 The first Gödel’s incompleteness theorem says that in every sufficiently  
120 rich<sup>1</sup> consistent first-order theory <sup>2</sup> there exist statements that are true<sup>3</sup>, but  
121 that cannot be proven within this theory. The second Gödel’s incompleteness  
122 theorem says that every sufficiently rich consistent first-order theory cannot  
123 prove its own consistency.

124 Human mathematicians – states the anti-mechanism argument based on

---

<sup>1</sup>By a “sufficiently rich” one means that the formal system is able to express arithmetic of addition and multiplication.

<sup>2</sup>A formal system, or a theory, is a collection of axioms together with rules of inference. The importance of using first-order logic is because of the completeness of this logic.

<sup>3</sup>A statement is true, when it is satisfied in the intended model of the theory.

125 Gödel's incompleteness theorems – can fruitfully work with Gödel's incom-  
126 pleteness theorems, thus those mathematicians use a the resources from the  
127 outside the theory, for instance, they are able to refer to the intended model  
128 of arithmetic or recognize that the human mind is consistent. Thus, human  
129 mathematicians outperform machines, because – unlike machines – they are  
130 able to include in their reasoning such external resources.

131 The intuition that humans could prove theorems which machines could  
132 not has been present already in Turing (1950 [47])<sup>4</sup> and in Post (1941 [31])<sup>5</sup>.  
133 One of the most famous voices exploring the anti-mechanism argument based  
134 on Gödel's incompleteness theorems against the computational theory of  
135 mind – next to Hofstadter (1979 [15]) and Nagel and Newman (1958 [24],  
136 1961 [25]) – Lucas (1961 [18], also 1968 [19], 1996 [21]) presented a “math-  
137 ematical proof” of man's superiority over a machine. Lucas extended the  
138 applicability of Gödel's incompleteness theorems from formal systems to hu-  
139 man subjects. On his view, humans are subjects to the same formal limits  
140 as machines. However, as Lucas observes, human mathematicians can prove  
141 Gödel's incompleteness theorem, which means, human mathematicians use  
142 extra-formal resources that enable them to perform such proofs.

143 Lucas' argument relies on the fact that Gödel's theorem(s) are formulated  
144 in purely formal terms. As Lucas observes himself, this is what differentiates  
145 the Gödel's results from the liar paradox. The liar paradox, stating that  
146 “This statement is untrue”, is “viciously self-referential, and we do not know  
147 what the statement is, which is alleged to be untrue, until it has been made,  
148 and we cannot make it until we know what it is that is being alleged to be  
149 false” (Lucas 1990 [20, page 2]). Unlike the liar paradox, Gödel's theorem  
150 is formulated within a full-blooded system where it is clearly defined, which  
151 sentences are true and what does it mean to be provable. Lucas' claims that  
152 the fact that an (idealised) human mind, even if cannot prove the Gödel's  
153 theorem(s) for the given theory, can – thanks to its additional non-mechanical  
154 skills – recognize this theorem as true in its system. In consequence, a human  
155 mind outperforms a machine.

156 Penrose in (1989 [26], 1994 [27]) extended Lucas' reasoning of a positive  
157 claim regarding the extra-formal resources available to humans that enable

---

<sup>4</sup>As reported by Krajewski, Turing believed that even if a machine cannot prove as much as humans can, it is still worth constructing robots.

<sup>5</sup>As reported by Krajewski, Post believed that man cannot construct a machine which can do all the things he can.

158 them to construct reasoning unavailable to machines. Penrose suggested that  
159 in the brain there exists physical basis of non-computable behavior, and he  
160 indicated quantum mechanics as a credible candidate. According to him  
161 quantum processes might explain not only reasoning of human mathemati-  
162 cians, but also consciousness.

163 A constructive criticism of Lucas-Penrose style argument got formulated  
164 by Putnam (1960 [32], Benacerraf (1967 [1]), Wang (1974 [53]), then later  
165 also by Boolos (1995 [2]) or Shapiro (1998 [43]). Penrose's version got criti-  
166 cised in particular by Feferman (1995 [11]), Putnam (1995 [34]) and Shapiro  
167 (2003 [44]). Krajewski claims that the ways of criticizing the Lucas-Penrose  
168 argument follow one of the two main lines [17, pages 5–6]:

- 169 • The mind is a machine and it is consistent, but it cannot prove the  
170 Gödel's sentence for itself.<sup>6</sup>
  
- 171 • The mind is a machine, but it is inconsistent, and Gödelian limitations  
172 do not apply to that.

173 Krajewski [17] refutes the Lucas-Putnam argument in yet another way:  
174 he observes that iterations of increasingly strong theories proving the cor-  
175 responding Gödel's sentences can be processed in a purely mechanical or  
176 computational manner available to both, humans and machines. In con-  
177 sequence, Krajewski claims that anti-mechanism is not implied by Gödel's  
178 incompleteness theorems stand alone. In addition, claims Krajewski, one  
179 needs to assume that humans have a privileged access to assessing consis-  
180 tency of the human mind. Krajewski claims that the argument fails because  
181 of the necessity of making this extra-formal assumption. This is so, because  
182 there is no formal way to account for the formal counterpart of assumptions.

183 Before I come back, in the last section, to Krajewski's rejection of the  
184 anti-mechanism argument, and my proposal how to shift the way of think-  
185 ing about the reasons for this rejection, I will now focus on the part which  
186 is particularly interesting for me, that is the *meta-theoretical* corollary to  
187 the anti-mechanism argument stating that humans cannot fully describe out  
188 notion of natural number.

---

<sup>6</sup>This line of argument comes already from Gödel, who distinguished *subjective arith-*  
*metic* that humans can do, and who believed that in *objective mathematics* full arithmetic  
is a consistent theory. He also believed that the concept of computation can be defined  
without referring to any domain of computation; these claims amount to Gödelian platon-  
ism. (Gödel \*1951 [13])

## 189 2. The concept of natural number and the concept of computation

190 I initiate my investigation into the nature of the extra-formal elements  
191 of the reasoning that enable the conclusion that the human mind is not  
192 computable, by discussing the corollary relating human inability to define  
193 the concept of natural number. Additionally, I extend the corollary of the  
194 claim that humans – for the similar reasons – cannot define the concept of  
195 computation. Finally, I present the view according to which the concept of  
196 natural number and the concept of computation are closely related.

197 The fact that every formal definition of the concept of natural num-  
198 ber leads to a necessary assumption from the outside of the formal system  
199 has been studied in the context of the view in philosophy of mathematics,  
200 called *structuralism*. According to structuralism, mathematics is “science of  
201 structures”, and while defining mathematical objects, one should first target  
202 their structural properties. For instance, while defining natural numbers, one  
203 should define the structure of natural numbers through relations they hold  
204 to each other, and not focus on individual properties of those elements.

205 Traditionally, structuralism defined natural numbers using the second-  
206 order Peano Arithmetic (PA2). PA2 is categorical and the class of (iso-  
207 morphic) models in which it is satisfied is identified with natural numbers.  
208 The usual way of criticising the use of second-order Peano Arithmetic to  
209 define natural numbers consists in saying that the underlying logic is “set  
210 theory in sheep’s clothing” (Quine 1970 [35, page 66]). Second-order logic  
211 has the ability, for instance, to express the information that two sets have  
212 the same cardinality. The concept of set is itself most frequently (implicitly)  
213 defined with a first-order axiomatic theory, such as  $ZF$ , that in its turn, is  
214 a subject of non-standard interpretations, Löwenheim-Skolem theorem, *etc.*,  
215 which makes its intended model “hidden” within a continuum of other non-  
216 intended models. Therefore, in order to define the concept of natural number  
217 with PA2, humans have two choices. They can get involved in a vicious circle  
218 of definitions, or an infinite regression of theorems, or they can use extra-  
219 formal resources and admit in an arbitrary manner that there is such a thing  
220 as an intended, (or a standard) model of set theory where exists the intended  
221 model of arithmetic.

222 Another, less known, version of structuralism, so called *computational*  
223 *structuralism*, proposes to distinguish the *standard* model of arithmetic from  
224 the continuum of non-standard models with the resources of PA1 only (Hal-  
225 bach & Horsten 2005 [14], Quinon & Zdanowski 2007 [36]). In order to



226 do that, defenders of computational structuralism suggest adding a meta-  
227 mathematical constraint regarding the computability of interpretation of  
228 functional symbols in the language, and then use Tennenbaum’s theorem  
229 in order to single out the standard model of arithmetic.

230 **Theorem 2.1 (Tennenbaum 1959).** *Let  $\mathcal{M} = \langle \mathbb{M}, +, \times, 0, 1, < \rangle$  be an enu-*  
231 *merable model of PA1, and not isomorphic with the standard model  $\mathcal{N} =$*   
232  *$\langle \mathbb{N}, +, \times, 0, 1, < \rangle$ . Then  $\mathcal{M}$  is not recursive.*

233 More explicitly why Tennenbaum’s theorem is relevant for the structural-  
234 ist way of thinking is visible in the transposition of the theorem:

235 **Theorem 2.2 (Tennenbaum transposition).** *Let  $M$  be an enumerable*  
236 *model of first-order Peano arithmetic. If the interpretation of addition and*  
237 *multiplication within  $M$  are computable then  $M$  is a standard model for arith-*  
238 *metic (a model with  $\omega$ -type ordering).*

239 One of the philosophically interesting consequences of the application of  
240 Tennenbaum’s theorem is that the set of models singled out with its help  
241 consists of those  $\omega$  models, where  $\omega$  is computable (Quinon & Zdanowski  
242 2007 [36]). Those models are called “intended” and form a proper subset of  
243 standard models.

244 The *intended* model of arithmetic<sup>7</sup>, is such a model where functions of  
245 addition and multiplication are interpreted as computable functions<sup>8</sup>. Ten-  
246 nenbaum’s theorem establishes a connection between a meta-mathematical  
247 property of being computable by arithmetical functions, and the order of  
248 the elements of the set of natural numbers. Thus, in the most general lines,  
249 computational structuralism is a position, according to which the concept of  
250 natural number and the concept of computation are closely related.

251 The usual way of criticising computational structuralism is, again, by  
252 pointing out at the vicious circle or infinite regression of definitions that  
253 threatens the proposed account of natural numbers. The criticism goes as fol-  
254 lows: in order to define the concept of natural number, one needs to use the  
255 concept of computation, whereas every concept of computation is defined on  
256 the domain of (some representation of) natural numbers. Thus, the vicious

---

<sup>7</sup>Intended models of arithmetic are identified up to a *computable* isomorphism.

<sup>8</sup>The model of arithmetic is intended for both theories PA1 and PA2

257 circle or the necessity to assume that there is an intended interpretation of  
258 what to compute means, or that the intended model of arithmetic is distin-  
259 guished from within other models.

260 Analogously, it is pretty straightforward that the concept of computation  
261 falls itself into a vicious circle, as in order to account for what “to compute”  
262 means, referring, for instance, “to be computed on a Turing Machine”, neces-  
263 sitates to account which entities are suitable for computing with (in the case  
264 of TM-computations, what can be the input for a Turing Machine). Since the  
265 question asked about the input precedes the definition of computing, which  
266 is just being given, one cannot use the concept of computing to define which  
267 sequences can be used for the input.

268 More precisely,

269 [...] the Church-Turing Thesis states that Turing Machines formally  
270 explicate the intuitive concept of computability. The description of Turing  
271 Machines requires description of the notation used for the INPUT and for the  
272 OUTPUT. The notation used by Turing in the original account and also notations  
273 used in contemporary handbooks of computability all belong to the most  
274 known, common, widespread notations, such as standard Arabic notation for  
275 natural numbers, binary encoding of natural numbers or stroke notation. The  
276 choice is arbitrary and left unjustified. In fact, providing such a justifi-  
277 cation and providing a general definition of notations, which are acceptable  
278 for the process of computations, causes problems. This is so, because the  
279 comprehensive definition states that such a notation or encoding has to  
280 be computable. Yet, using the concept of computability in a definition of  
281 a notation, which will be further used in a definition of the concept of  
282 computability yields an obvious vicious circle. (Quinon 2018 [page 338][38]).  
283  
284  
285

286 In this section, I explained similarities between the process of defining  
287 the concept of natural number, the process of accounting for the concept of  
288 computation, and the formulation of an anti-mechanism argument based on  
289 Gödel’s incompleteness theorems. All those contexts are related because the  
290 way out of the definitional vicious circles proper to the definitional processes  
291 within formal theories, through the necessity of assuming an additional non-  
292 formal, meta-theoretical knowledge. In the next section, I will develop on  
293 the phenomena of vicious circles and regression ad infinitum.

294 **3. Nested vicious circles**

295 Quinon (2018 [38]) proposes a taxonomy of what can be called “deviant  
296 encodings”, that is those encodings – or in different words, sequences of  
297 symbolic representations of natural numbers – which are non-computable,  
298 but which are formally indistinguishable from computable encodings. For  
299 instance, in its simplest form the problem presents itself as follows:

300 The problem in its purely syntactical version can be formulated  
301 as follows. In a definition of Turing computability, one of the  
302 aspects that needs to be clarified is the characterization of nota-  
303 tion that can be used as an input for a machine to process. If  
304 a Turing Machine is supposed to explicate the intuitive concept  
305 of computability it is necessary to explain, which sequence of num-  
306 erals can be used as an input without the use of the concept  
307 of computability. That means, we cannot simply say: “sequences  
308 that can be used as input are the computable ones” as we have  
309 not yet defined what means “to be computable”. (Quinon 2018  
310 [38, page 340]).

311 Deviations refer to non-computable sequences that cannot be distinguished  
312 within the general formal context from sequences that are computable and  
313 can be used in computations. In this paper, I use the expression “deviant  
314 encoding” independently of the ontological framework within which natural  
315 numbers are understood. Quinon (2018 [38]) claims that the phenomenon of  
316 deviant encodings persist independently of which ontological status we as-  
317 sign to objects of computations (*e.g.*, natural numbers, sequences of symbols,  
318 *etc.*). Quinon (2018 [38]) hypothesizes that the phenomenon of deviant encod-  
319 ings persists independently of the philosophical standpoint and provides an  
320 analysis of following simplified standpoints: (i) purely mechanical/syntactical  
321 approach (nominalism, entwined mathematical concepts); (ii) notations have  
322 meanings (mild realism); (iii) semantics comes first (radical realism, platonic  
323 insight).

324 The study of conceptual “deviations” is conducted for a simplified frame-  
325 work where:

- 326 • on the syntactic level there are uninterpreted inscriptions, and where  
327 functions are string-theoretical generating strings values from string  
328 arguments;

- 329     • on the semantic level there are interpretations that can range from  
330       the conceptual content ascribed to initially uninterpreted symbols, to  
331       Platonic abstract objects, and where functions are number-theoretical  
332       sending numbers to numbers;
- 333     • between the two levels there is defined a function of denotation.

334 Deviations occur on each level. Thus, there exist “deviant encodings” devi-  
335 ations that happens on the syntactic level; “deviant semantics” deviations  
336 that happens on the semantic level; “unacceptable denotation function” devi-  
337 ations of the denotation function.

338     The simplified framework is inspired by Shapiro (1982 [42]), who dis-  
339 tinguishes string-theoretic functions from number-theoretic functions and  
340 searches for “acceptable”, that is “non-deviant”, ways of associating their  
341 domains. The framework is further used by other researchers. Rescorla (2007  
342 [41]) uses it to study behaviour of denotation functions which associate nu-  
343 merals (symbolic representations of natural numbers) to natural numbers  
344 (abstract entities) in a non-computable manner. There is a continuum of  
345 such mappings.

346     Differently use the expression “deviant encodings” Copeland and Proud-  
347 foot (2010 [7]) for whom the deviations relate encodings, or enumerations, of  
348 Turing Machines. The authors claim that a deviant encoding happens when  
349 the omniscient programmer “winks at us” to let us know when the number  
350 of a Turing Machine (from some standard encoding of Turing Machines),  
351 which is being currently processed by some sort of the Halting Machine (a  
352 machine computing which Turing Machines stop on an input 0), refers to a  
353 machine that stops. In this way, the Halting Machine computes the halting  
354 function, which is an uncomputable function. The “wink” of the omniscient  
355 programmer gets encoded in the syntactic structure of the numerals: the nu-  
356 merals representing machines that stop, have a special form, – for instance –  
357 are even (their general syntactical form can be reduced to “2n” where “n” is  
358 any numeral). Copeland and Proudfoot mean by a deviant encoding such a  
359 standard enumeration of Turing Machines where the encoding is enriched by  
360 an extra-formal feature impersonated by the omniscient programmer (a Tur-  
361 ing oracle). This is a specific case of a more general problem where deviant  
362 encodings refer to encodings representing natural numbers.

363     An occurrence of the phenomenon of deviant encodings involving all the  
364 levels, is the case of the Semantical Halting Problem (van Heuveln 2000  
365 [52]). Imagine, you have encoded Turing machines with some standard –

366 computable, thus non-deviant – encoding, and that you believe that symbols  
367 have meanings or interpretations. It can happen that even if your syntax is  
368 generated in a recursive manner, your semantics is not following any recursive  
369 rules. The Halting Machine that processes encodings of Turing Machines is  
370 designed to process information on syntax in an algorithmic manner. If  
371 inputted with a given non-computable enumeration of Turing machines, the  
372 machine will process those non-computable encodings as if it were a standard  
373 notation. Again, there is no effective way of defining which semantics are  
374 acceptable and which are deviant.

375 I call “nested vicious circles” the hierarchies of vicious circles that keep  
376 reappearing at every stage of syntactical and semantic complexity of the pre-  
377 sented picture.

378 To give an example of a philosophical position outside the strict theoret-  
379 ical context discussed in this paper, the phenomenon of deviant encodings  
380 appears as well in the case of concrete computations.

381 In our ordinary discourse, we distinguish between physical sys-  
382 tems that perform computations, such as computers and calcul-  
383 ators, and physical systems that don’t, such as rocks. Among  
384 computing devices, we distinguish between more and less power-  
385 ful ones. These distinctions affect our behaviour: if a device is  
386 computationally more powerful than another, we pay more money  
387 for it. What grounds these distinctions? What is the principled  
388 difference, if there is one, between a rock and a calculator, or be-  
389 tween a calculator and a computer? Answering these questions  
390 is more difficult than it may seem. (Piccinini 2010 [28])<sup>9</sup>.

391 In (2020 [40]), Quinon notes that the phenomenon of nested vicious cir-  
392 cles, relating the concept of computability, do not disappear in the case of  
393 explicit inter-definability between the concept of natural number and the  
394 concept of computation, as established by computational structuralism. As  
395 I already described it above, the criticism of computational structuralism  
396 consists in pointing at the choice between the definitional vicious circles or  
397 the necessity of making extra-formal arbitrary assumptions.

398 The way of extra-formal assumptions is investigated by Button and Smith  
399 (2012 [3]) who observed that when the concept “natural number” is expli-

---

<sup>9</sup>See also Piccinini (2015 [29]).

400 cated for, the concepts used in this explication, such as “to compute” or  
401 “finite” need to be accounted for on their turn, *etc.* In consequence, claim  
402 the authors, this problem cannot be tackled by offering more mathematics.  
403 An arbitrary decision regarding the meaning of some concept is necessary for  
404 the argument from the Tennenbaum’s theorem to work. However, as they  
405 claim in a slightly undermining way, this is a philosophical problem:

406        Suffice it to note that our discussion of Tennenbaum’s Theorem il-  
407        lustrates a familiar moral: philosophical problems which are sup-  
408        posedly generated by mathematical results can rarely be tackled  
409        by offering more mathematics. (Button & Smith 2012 [3, page  
410        120]).

411        Dean (2014 [9]) is similarly sceptic when it comes to the purposefulness  
412        of using Tennenbaum’s theorem to formally single out the standard model of  
413        arithmetic. However, differently to Button and Smith, Dean develops a full  
414        fledged philosophical position. It is a Putnam-style model-theoretic realism  
415        for the concept of computation (see Putnam 1980 [33]). Dean claims that  
416        there is no point in trying to find external arguments to distinguish between  
417        various standard and non-standard models neither of arithmetic, nor of recur-  
418        sive theory. We should rather use the richness of the model-theoretic universe  
419        for studying structural properties of the concept of computation. Dean claims  
420        that Tennenbaum’s phenomenon shows that there exists continuum of pairs:  
421        model of arithmetic and computation in this model of arithmetic. In conse-  
422        quence, the Tennenbaum’s result instead of contributing to singling out the  
423        standard model of arithmetic, it indicates that there exist non-computable  
424        *omega*-models of arithmetic (the so called deviant or weird permutations)  
425        with a corresponding concept of computation defined within the model.

426        The vicious circle faced by computational structuralism, differs from the  
427        vicious circles that are the focus of Quinon (2018 [38]). There, I was only  
428        concerned by the concept of natural number being indirectly involved in  
429        the definition of what “to compute” means. Conceptual structuralism needs  
430        to handle a slightly more elaborate idea. Its objective is to explicate the  
431        concept of natural number, identified with the standard model of arithmetic.  
432        Its solution consists in using the idea that natural numbers, and in particular  
433        those which are defined by Peano’s axioms, are the entities used for counting  
434        and computing. In consequence, natural numbers are defined in terms of  
435        computations. However, and this is where the vicious circle arises: one of  
436        the characteristic features of the concept of computation is that computation

437 is *always* defined on some given domain.<sup>10</sup> This domain is always identifiable  
438 with the structure of natural numbers. I discuss the nested vicious circles in  
439 this context in Quinon (2020 [40]).

#### 440 4. Conceptual engineering and conceptual fixed points

441 One of the promising ways out of the impasse consists in embracing that  
442 the circularity in the account of what “to be computable” and what “natural  
443 number” mean is due to limitations of conceptual analysis. Similarly to other  
444 scientific concepts, when analysis is conducted within the strict scope of a  
445 given formal theory, one often ends up with a necessity to use the concept  
446 which is being defined in the account of some concept used for its definition.  
447 Philosophers and logicians see in this feature of conceptual analysis both, an  
448 advantage that enables us to understand more about the conceptual structure  
449 of the world (see Dean (2014 [9]), and a problem that blocks science from  
450 progress (see Maddy (2007 [22])). Rescorla (2007 [41]) identifies problems  
451 with conceptual analysis related to the concept of computation. In their  
452 paper [3] Button and Smith claim that Tennenbaum’s theorem is of no use  
453 for a philosopher who wants to distinguish the standard model from other  
454 possible models of arithmetic.

455 Quinon (2018 [38]) suggests that there is no fully satisfactory way out  
456 from vicious circles in definitions, resulting from conceptual analysis. Ap-  
457 proaching the concept of computation and the concept of natural number  
458 from another methodological perspective, seems to be more fruitful. For  
459 instance, in the recent years a particular type of conceptual work gained  
460 quite a bit of popularity, it is called *conceptual engineering*. What I try to  
461 convey in this section is that the new research on conceptual engineering  
462 actually provide additional insight into the possible ways of thinking about  
463 the non-mathematical or non-formal knowledge.

464 According Cappelen (2018 [4]), conceptual engineering is concerned with  
465 the assessment and improvement of concepts. As highlight Cappelen and  
466 Plunkett (2020 [30, page 3]):

467 “since it’s unclear and controversial what concepts are (and whether  
468 there are any), it’s better to broaden the scope along the following

---

<sup>10</sup>A non-realised Gödel’s objective consisted in finding an “absolute” concept of computation, *i.e.*, such a concept of computation that does not depend on any domain.

469

lines:

470

**Conceptual Engineering** = (i) The assessment of representational devices, (ii) reflections on and proposal for how to improve representational devices, and (iii) efforts to implement the proposed improvements.”

471

472

473

474

Researchers involved in developing the methodology of conceptual engineering realised that the method reaches its limits when concepts which are fundamental to the given theory are being scrutinised. They call it “conceptual fixed points”. The most extensive reflection has been done in the area of ethics (Cappelen et al. 2020 [5]), but Eklund (2015 [10]) extends it to formal contexts and concepts such as “truth”, “belief”, or “existence”. In addition to traditional arguments used in the ethical contexts, such as “Kantian philosophy [with its regulative ideas], or from a naturalistic philosophy according to which what is innate severely constrains which concepts we can use”, Eklund considers basic formal concepts in the spirit of rigid designators.

485

In moral philosophy, “the moral fixed points” are those moral propositions that are moral truths that always need to be incorporated in a moral system. A normative system which fails to incorporate such propositions is not a moral system, but a normative system of some other kind. The flag example of such a moral fixed point is a proposition “It is wrong to engage in the recreational slaughter of a fellow person” (Cueno & Shafer-Landau 2014 [8]).

491

Eklund (*e.g.*, 2015 [10, chapter 5]) extends this phenomenon to frameworks outside moral philosophy and, as he calls it, the “thinnest” normative words like “good”, “right”, “ought”. Eklund observes that in each conceptual framework, there exist concepts that are difficult, if not impossible to engineer. “Truth” is one of those concepts. People care about truth, writes Eklund, and they do not care about some conceptually engineered concept “truth\*”. In consequence, truth is a concept that should keep a fixed position in a conceptual framework, and refer to the natural kin of assertions and beliefs. Similarly, “existence” is a conceptual fixed point. Eklund opposes the claim from in the contemporary meta-ontological debate, where it is assumed “that there are alternative notions of existence that can be employed”. He claims that, similarly as in the case of “truth”, a conceptual framework that would result from adapting a conceptually engineered concept of “existence” would need to adjust its other key concepts in such a way that the resulting

504



505 framework would be isomorphic to the initial one. Thus, “One cannot, so to  
506 speak, *selectively* engineer the quantifier”.

507       Suppose we set out to conceptually engineer truth. Insofar as the  
508       job description of truth is that of being the property our beliefs  
509       and assertions aim at, the engineering project would be that of  
510       finding a property more adequate to that job description. But  
511       by what has been noted about Stich’s argument, it is hard even  
512       properly to conceive of a practice of belief or assertion that is  
513       guided by a different property. (Eklund 2015 [10, page 378])

514       There is one last thing that I consider worth mentioning while talking  
515       about conceptual fixed points and mathematical concepts, in particular the  
516       concept of computation, that is a possible proximity between conceptual fixed  
517       points and fixed points that are traditionally analysed in mathematics in the  
518       context of diagonalisation. At the first sight, they do not have much in com-  
519       mon<sup>11</sup> as conceptual fixed points relate mostly to the cross-model intended in-  
520       terpretation of a concept, whereas diagonalisation is about self-reference and  
521       vicious circles. Conceptual fixed points are concepts interpreted in, what we  
522       call in philosophy of mathematics, their intended models. In different words,  
523       a fixed point consists of the pair *the engineered concept* corresponding to  
524       the intended meaning of the concept, or – to borrow Eklund’s expression –  
525       the interpretation that “people care about”, and *a possible world of inter-*  
526       *pretation*, which actually corresponds to the intended model of this concept.  
527       Both, the concept of natural number and the concept of computation are in  
528       this sense conceptual fixed points. A more careful look should be applied to  
529       those two phenomena, but in this paper I will just leave it without further  
530       comment<sup>12</sup>.

## 531 **5. The Church-Turing Thesis as a Carnapian explication**

532       Another methodological framework that offers a solution for conceptual  
533       structure escaping conceptual analysis is the method of Carnapian explica-  
534       tion. Quinon (2019 [39]) explores the idea that the structure of the concept

---

<sup>11</sup>I might be wrong, but I will not try to sort it out in this paper.

<sup>12</sup>If you want to get a more formal description of this phenomenon, you can think of hybrid modal logics which provide a framework for thinking of epistemic access to other possible worlds from the perspective of the selected distinguished world.

535 of computation, accounted for with the Church-Turing thesis, is the best un-  
536 derstood through the method of explication. This section is devoted to the  
537 presentation of the method of explication for the concept of computation,  
538 and also for the concept of natural number.

539 Treating the concept of computation, as accounted for in the Church-  
540 Turing thesis, as a Carnapian explication has multiple advantages, namely,  
541 it overcomes problems of conceptual analysis; it explains how one intuitive  
542 concept of what “to be computable” means can be translated into a multi-  
543 tude of extensionally equivalent formal concepts of “to be computable” in  
544 a specific formal concept means; it finally provides a ground for thinking of  
545 mathematical or formal concepts as “open-textures” evolving through the  
546 time (Makovec & Shapiro 2019 [23]); it also relates the initial intuitive pre-  
547 scientific concept with the formal concept, because an explication relies on  
548 an existing meaning, and offers a specification which suites best possible in  
549 a given context.

550 An explication in the Carnapian sense consists in introducing new formal  
551 concepts to the scientific language coined on the basis of everyday concepts.  
552 In different words, it is a procedure of transformation of an inexact pre-  
553 scientific concept into a scientific one. Moreover, an explication consists  
554 in providing a scientific concept within a given context, within an existing  
555 theory. It is done in two steps:

- 556 • The clarification of the explicatum
- 557 • The specification of the explicatum

558 The rationale for clarification is that a given term may have many different  
559 meanings in ordinary language. Unless one of these meanings is clearly picked  
560 out from the start and the context of its use is clearly indicated, it is unlikely  
561 that the method of explication will yield a useful result. Clarification serves  
562 this purpose. As Carnap explains,

563 [a]lthough the explicandum cannot be given in exact terms, it  
564 should be made as clear as possible by informal explanations and  
565 examples. (Carnap 1950 [6, page 3]).

566 Quinon (2019 [39]) highlights the importance of the clarification stage, the  
567 stage which has traditionally been underestimated.

568 A clarification of the explicandum enables the next step of the explica-  
569 tion process, a specification of the explicatum and formulation of the exact  
570 concept in the targeted context.

571 Since most often several clarifications can be foreseen, and several sci-  
572 entific contexts are available, one pre-scientific concept can be explicated in  
573 various manners. In order to decide which explication is the most successful,  
574 Carnap proposes four criteria that can be applied for assessing the value of  
575 an explication, and also for comparison between available options.

- 576 • *Similarity to the explicandum* most of the cases in which the explican-  
577 dum has so far been used, the explicatum can be used; however, close  
578 similarity is not required, and considerable differences are permitted.
- 579 • *Exactness* the rules of use of explicatum have to be given explicitly  
580 and precisely, for example, by providing a concept with the formal  
581 definition.
- 582 • *Fruitfulness* shall be “useful for the formulation of many universal state-  
583 ments”.
- 584 • *Simplicity* an explication should be as simple as the previous three  
585 allow it.

586 I think that it is worth investigating whether abandoning the path of  
587 analysis and taking the path of explications could offer an additional insight  
588 into the conceptual structure of formal concepts, and also informal concepts  
589 lying in the foundations of their formalization. The idea that every formal  
590 concept is, – at least in subjective arithmetic (to borrow Gödelian termi-  
591 nology) – grounded upon, or issued from, everyday intuitive, pre-scientific  
592 concept. The next section, is devoted to a preliminary investigation into a  
593 possibility of extending the idea that the method of explication, consisting  
594 in building up the formal concept out of the intuitive concept, is anyhow  
595 relevant to the anti-mechanism argument against the computability of mind  
596 using Gödel’s incompleteness theorems.

597 Both intended interpretations determined in consequences of accepting  
598 conceptual fixed points solution and the choice of the formal aspect and  
599 the formal context at the stage of the concept clarification in the process  
600 of Carnapian explication, share a similar threat. In the case of fixed point  
601 solution and in the case of clarification an agent needs to take an arbitrary  
602 decision regarding the intended interpretation.

## 603 **6. Theory of mind and computations**

604 In this section, I propose an additional complication to the method of  
605 Carnapian explication, which is a temporary, or a phylogenic, aspect of con-  
606 ceptual development.

607 The method of Carnapian explication enables introducing new formal  
608 concepts to the language by transforming an intuitive pre-scientific concept  
609 into a new scientific concept within some formal concept. Usually, at the  
610 stage of clarification one chooses the formal meaning that the formalisation  
611 of the intuitive pre-scientific concept and also the targeted formal concept.  
612 What I propose in this section is an additional dimension to the clarification  
613 stage: a relativisation to the phylogeny of the targeted concept. At the stage  
614 of clarification, in addition to deciding which aspect of the intuitive concept  
615 one wants to formalise one needs to realise that each concept develops. The  
616 phylogenic development of the concept of natural number and the concept  
617 of computation is studied in Shapiro on open texture (2013 [45]).

618 The relation between the concept of computation and the concept of  
619 natural number underwent a very dynamic development. In consequence,  
620 grew the set potential clarifications of intuitive concepts of computation and  
621 of natural numbers. What is interesting from my perspective is that com-  
622 putability is today an expected feature of natural numbers. Natural numbers  
623 are those mathematical entities that are all days long used for enumerating  
624 and computing, for programming, and in various sorts of logistic projects as  
625 an underlying discrete structure. Both concepts get increasingly important  
626 in the everyday life of our society. This is called digitalisation.

627 Various areas of digitalisation are additionally reinforced by the fact that  
628 computationalism – even if its formal details are still discussed by philoso-  
629 phers, mathematicians and logicians – is today the mainstream theory of  
630 mind. This process is described by Turkle (1984/2004, 2011, 2015) who  
631 studies how concepts from computer sciences and robotics get into common  
632 language and how they change ordinary people’s approach to inter-personal  
633 relations or ethical questions.

634 According to Turkle the intensity in which digitalisation of the everyday  
635 life develops is strongly connected to the fact that computational language  
636 was first used to reformulate our perception of our own mind and our con-  
637 sciousness.<sup>13</sup>

---

<sup>13</sup>Turkle’s earlier work related similar development of conceptual trends in explanation

638 When Turkle speaks about her experience with the digitalised society, she  
639 compares the two experiences:

640 My experience at MIT impressed me with the fact that some-  
641 thing analogous to the development of a psychoanalytic culture  
642 was going on in the worlds around computation. At MIT I heard  
643 computational metaphors used to think about politics, education,  
644 social process, and, most central to the analogy with psychoanal-  
645 ysis, about the self. (Turkle 1984 [49, page 305])

646 She sees in it a first step in the cultural assimilation of a new way of  
647 thinking:

648 The essential question in such work is how ideas developed in the  
649 world of high science are appropriated by the culture at large. In  
650 the case of psychoanalysis, how do Freudian ideas move out to  
651 touch the lives of people who have never visited a psychoanalyst,  
652 people who are not even particularly interested in psychoanalysis  
653 as a theory? In the study of the nascent computer culture, the  
654 essential question was the same: how were computational ideas  
655 moving out into everyday life? (Turkle 1984 [49, page 305]).

656 She searches how “the idea of mind as a program enters into peoples sense  
657 of who is the actor when they act”. A model of mind that is adapted by the  
658 society influence how people think about their frustrations and disappoint-  
659 ments, their relationships with their families and with their work [23, page

---

of phenomena of everyday life that had place in France in the 1960s and 1970s in consequence of the spread of psychoanalytical ideas, see her book “Psychoanalytic Politics: Jacques Lacan and Freud’s French Revolution” from 1978 [48]). In “The Second Self: Computers and the Human Spirit” (1984/2005 [49]) Turkle describes these changes that gets into general culture from the digitalisation and robotics in the same way as “psychoanalytic culture” penetrated structures of the general social and political life in France: “Psychoanalytic language spread into the rhetoric of political parties, into training programs for schoolteachers, into advice-to-the-lovelorn columns. I became fascinated with how people were picking up and trying on this new language for thinking about the self. I had gone to France to study the psychoanalytic community and how it had reinvented Freud for the French taste, but I was there at a time when it was possible to watch a small psychoanalytic community grow into a larger psychoanalytic culture.” (Turkle 1984 [49, pages 304-305]).

660 305]. On the other hand, says Turkle, computer became a new constructed  
661 object - “a cultural object that different people and groups of people can  
662 apprehend with very different descriptions and invest with very different at-  
663 tributes. Ideas about computers become easily charged with personal and  
664 cultural meanings” (Turkle 1984 [49, page 308]).

665 In her other books, Turkle studies human attachment to objects. In the  
666 volume of essays “Evocative Objects: Things We Think With” she speaks  
667 about attachment that people, many of her friends, developed with physical  
668 objects. In her book, “Alone Together” Turkle (2011 [50]) extends her obser-  
669 vations to different types of automated artificial agents, such as virtual agents  
670 mediated by electronic support, or robots. In the series of social experiments,  
671 where she asks her subjects to interact with an automated artificial agent,  
672 she observes that the stronger attachment develops in the most vulnerable  
673 members of our society, such as neglected children with unfulfilled emotional  
674 needs, or as old people suffering from the lack of human interactions. Our  
675 natural inclination to form emotional attachment with humans, and with  
676 objects in the absence of humans, might soon lead to even more human-AI  
677 interactions. Those interactions are obviously structured in a very particu-  
678 lar, very automated, way, which even more strongly influence digitalisation  
679 of the language we use.

680 Krajewski makes a similar observation in the last section of the paper.

681 Our attitude toward the arguments of Lucas, Penrose, and others  
682 is shaped mostly by our general vision of machines and minds.  
683 And this vision adjusts with changes of civilization. For the youth  
684 of today, if I may judge from listening to my students, our com-  
685 puterized world makes it easier to accept the idea that anything  
686 is mechanizable – including the mind. [17, page 32]

687 I propose a hypothesis that at least a part of the confusion regarding  
688 the specificity of the conceptual structure of the concept of computation  
689 contributes to the confusion regarding nature of human reasoning and the  
690 human mind. In consequence, I claim that at least partially the “feeling”  
691 that there exist non-computational processes come from the complexity of  
692 the conceptual structure of the concept of computation.

693 **7. The Lucas-Penrose argument and extra-formal concepts**

694 Let me now come back to the anti-mechanism argument against com-  
695 putability of mind based on Gödel’s incompleteness theorems. In this sec-  
696 tion, I reconstruct Krajewski’s argument that in order to render the anti-  
697 mechanism argument work, one needs to add an extra-formal assumption  
698 regarding the consistency of the underlying theory, that is, the theory cor-  
699 responding to the human mind. I present the core of Krajewski’s criticism:  
700 it is not possible to formalize the extra-formal assumption and hence, the  
701 whole Lucas’ argument is fallacious. I disagree with Krajewski’s claim that  
702 a formalization of the extra formal assumptions is not possible. There are  
703 contemporary philosophical methods that might enable formulation of such  
704 a formalization. In previous sections, I shown how the methodological and  
705 conceptual framework issued from Carnapian explications. Instead, I focus  
706 on another problem which issues form an internal characteristic of formal  
707 contexts, namely on this aspect of the argument which leads to a circular  
708 reasoning in the reasoning, *i.e.*, in order to show that the human mind ( $T_{HM}$ )  
709 outperforms a machine ( $T_M$ ) one needs to assume that the human mind is  
710 consistent and knows it (and in this way outperforms a machine that can  
711 never “know” if it is consistent or not). Again, I already discussed how the  
712 method of conceptual engineering enable structuring thinking of extra-formal  
713 assumption and the resulting circular reasoning.

714 In the second part of this section, I continue my investigation of possible  
715 extra-formal assumptions relative to the anti-mechanism argument based on  
716 Godel’s incompleteness theorems.

717 The Lucas’ anti-mechanism argument based on Gödel’s incompleteness  
718 theorems consists of two parts. Firstly, Gödel’s results establish that each  
719 sufficiently rich consistent theory admits a Gödel’s sentence and also that  
720 none such theory can prove its own consistency.

Let  $T$  be a consistent theory containing arithmetic, let  $\varphi_T$  be the Gödel’s sentence for the theory  $T$ .

$$Con(T) \Rightarrow T \not\vdash \varphi_T$$

$$Con(T) \Rightarrow T \not\vdash Con(T)$$

721 Moreover, it is broadly known that an inconsistent theory proves any sen-  
722 tence, but Gödel’s incompleteness theorems do not apply to an inconsistent  
723 theory.

Secondly, a human mathematicians can work with subsequent increasingly stronger theories,

$$\begin{aligned}T_1 &= T \cup \text{Con}(T) \\T_2 &= T_1 \cup \text{Con}(T) \\&\vdots \\T_{n+1} &= T \cup \text{Con}(T)\end{aligned}$$

724 which – for some defenders of an anti-mechanism argument – signifies that  
725 human mathematicians outperform machines. Krajewski objects this view  
726 claiming that the construction of the hierarchy can be fully mechanised. In  
727 consequence, he claims that the ability to construct and work with the hier-  
728 archy of increasingly stronger theories alone is not sufficient for formulating  
729 the anti-mechanism argument. As stated by Krajewski [17, page 30–31],  
730 additional assumptions are missing.

731 In addition to Gödel’s results, at least two assumptions that are  
732 not self-evident are used in the above reasoning. First, every  
733 exact proof of our consistency can be formalized, second, it is  
734 possible to express “our consistency.” [...] If this is accepted,  
735 one could question the second point. It is not clear at all how one  
736 can express “our consistency.” Basically there are two options  
737 to express this: either (i) by the common sense statement “I am  
738 consistent” or (ii) by a formal counterpart to this statement. Let  
739 us consider them in turn.

740 In case (i) we refer to a common sense statement, which have no  
741 connection to formal considerations. Hao Wang [in 1974, pages  
742 317-320] reflected on just this statement and believed that it is  
743 not provable. [...] If that were possible, it would mean that we  
744 are not machines, or that we are not even equivalent to machines  
745 in the realm of proof-producing reasoning. We certainly may  
746 believe that, but it is no more than a general feeling.

747 In case (ii) we consider the formal counterpart to a loose state-  
748 ment expressing consistency; [...] The usual meaning of the  
749 statement refers to the will to avoid contradictions, to the re-  
750 liability of our vision of the world, to the claim that the methods  
751 used by mathematicians are unailing. The sentence *Cons* or any



752 other similar arithmetical formula is rather far from those ideas.  
753 Thus, while something is strictly proved, it is unclear to what  
754 extent the conclusion conveys our consistency.

755 Krajewski's reasoning can be reconstructed as follows. Applying the for-  
756 mal predicate "being consistent" can only apply to a formal theory. Applying  
757 the formal predicate "being consistent" to anything else than a formal theory  
758 is a categorical mistake. In consequence, if consistency is to be predicated of  
759 the human mind, the mind must have certain formal properties and needs to  
760 be identified with a theory. There exist the following options:

- 761 • If human mind is a theory and it is consistent, then as to all other the-  
762 ories, a Gödel's sentence applies to it and the human mind encounters  
763 the same constraints as any theory (a machine).
- 764 • If the human mind is a theory and it is inconsistent, then Gödelian  
765 argument limitations do not apply at all.

766 If the human mind is a theory, a human disposing of a mind cannot know  
767 – from the formal point of view – if it is consistent or not. In consequence,  
768 in order to prove that the human mind outperforms a machine, a second  
769 extra-formal additional assumption needs to be made. It has to be assumed  
770 that the human mind is indeed consistent. This assumption can be done in  
771 one of the two ways. "Case (i)", "I am consistent" cannot be formalised.  
772 "Case (ii)", there exists a formal counterpart of "I am consistent".

773 My analysis of "case (i)" is in line with the analysis of Krajewski. If "I  
774 am consistent" is an informal statement, it is useless for any formal proof.  
775 And here we speak of being able to *prove* more than a machine. Whereas  
776 Lucas' argument is supposed to be a formal proof of the superiority of the  
777 human mind over a machine.

778 My analysis of "case (ii)" differs from Krajewski's analysis. His argument  
779 resumes to the idea that each formalisation of the informal "I am consistent"  
780 remains – maybe more informed or more precise – but still an informal ac-  
781 count. As such it is useless for any formal proof. I think that conclusion from  
782 (ii) is different. An agent *can* find a formal counterpart of the statement "I  
783 am consistent", or rather "the theory constituting my mind is consistent".  
784 The framework of Carnapian explications enables us to understand how it  
785 can be done.

786 I also assume that an agent *can* recognise her own consistency. This  
787 insight is available to a human being, while it is – on the ground of the second  
788 Gödel’s incompleteness theorem – unavailable to a machine. This extra-  
789 formal assumption is necessary for formulating an anti-mechanism argument  
790 against computability of mind. This is also exactly at this point where a  
791 vicious circle occurs. We are just being in the act of proving that the human  
792 mind outperforms a machine, and here we are assuming exactly this.

793 Another possible extra-formal assumption that can be made to enable the  
794 anti-mechanism argument based on Gödel’s incompleteness theorem, is the  
795 ability to refer to the intended model of arithmetic<sup>14</sup>. Instead of assuming  
796 that the human mind is consistent, that means, assuming that the theory  
797 underlying all human reasoning is a consistent theory which does not prove  
798 both a  $\varphi$  and a  $\neg\varphi$ , for every  $\varphi$ , in order to use Gödel’s incompleteness theo-  
799 rems to support anti-mechanism argument, one can assume that the human  
800 mind is able to refer to the intended model of arithmetic. The assumption  
801 that the human mind can refer to the intended model of arithmetic disables  
802 the possibility that the Gödel sentences get to have non-standard Gödel’s  
803 numerals.

804 In the way it is usually interpreted – in particular in the context of philo-  
805 sophical argumentation supporting anti-mechanism argument that the hu-  
806 man mind is non-computable – Gödel’s incompleteness theorems provide  
807 us with the information from the perspective of a formal system. The se-  
808 mantical aspect is taken for granted. When the model theoretical reason-  
809 ing is applied, Gödel’s incompleteness theorems indicates that there exist  
810 non-standard models where the (non-standard) Gödel number of the proof  
811 of Gödel’s incompleteness theorems have its (semantical) reference. It also  
812 means that there exist models where the Gödel (non-standard) number of  
813 the proof of the negation of Gödel’s first theorem, has an interpretation as a  
814 (non-standard) natural number.

815 What is famously referred to by Gödel’s platonism is his belief that there  
816 is a model of arithmetic in which all arithmetical truths are satisfied. This is  
817 obviously not the intended model of arithmetic that humans have a privileged  
818 cognitive access to, but the model of arithmetic in the objective mathematics

---

<sup>14</sup>The intended model is intended for both *PA1* and *PA2* and for this reason I do not make a distinction between the intended model of *PA1* and the intended model of *PA2*. I can think of a philosophical position that makes such a distinction, but for my purpose that would unnecessarily complicate my presentation.

819 (Gödel \*1951 [13]).

## 820 **Conclusions**

821 Additional to the critical analysis of Krajewski’s rejection of the anti-  
822 mechanism based on Gödel’s incompleteness theorems to which I suggest  
823 some possible improvements, my paper is sympathetic to the idea that cer-  
824 tain key concepts in formal contexts naturally fall into circular or infinite  
825 reasonings. In this way, I try to shift attention from the theory of the hu-  
826 man mind and consciousness, to the study of the conceptual structure of the  
827 language.

828 In my paper, I explored similarities between various formal contexts in  
829 which key concepts fall into a vicious circle of reasoning. I looked at the  
830 formalisation of the concept of natural number, of the concept of compu-  
831 tation, and at the concept of consistency in the context of Gödel’s incom-  
832 pleteness theorems. I suggested that the way to switch from an informal  
833 pre-scientific concept to a full-blooded formal scientific concept formulated  
834 in an adequate formal context is best modeled by Carnapian explications.  
835 I also suggested that the phenomenon of conceptual fixed points offers a  
836 methodological framework to think of intended interpretations necessary to  
837 jump out of circularity.

## 838 **References**

- 839 [1] Benacerraf, P. (1967). God, the Devil, and Gödel. *Monist* 51: 9-32.
- 840 [2] Boolos, G. (1995). Introductory note to \*1951. In: Gödel, K., *Collected*  
841 *Works, Volume III, Unpublished essays and lectures*, Feferman S., et al.  
842 (eds.), Oxford University Press (1995): 290—304.
- 843 [3] Button, T., Smith, P. (2012): The Philosophical Significance of Tennen-  
844 baum’s Theorem. *Philosophia Mathematica* 20(1): 114–121.
- 845 [4] Cappelen, H. (2018). *Fixing Language: An Essay on Conceptual Engi-*  
846 *neering*. Oxford University Press.
- 847 [5] Cappelen H., Plunkett D. & Burgess A. (eds.) (2020), *Conceptual Engi-*  
848 *neering and Conceptual Ethics*. Oxford University Press.

- 849 [6] Carnap, Rudolf (1950), *Logical Foundations of Probability*. Routledge  
850 and Kegan Paul.
- 851 [7] Copeland, J. & Proudfoot, D. (2010). Deviant encodings and Turing's  
852 analysis of computability. *Studies in History and Philosophy of Science*  
853 41: 247–252.
- 854 [8] Cuneo T. & Shafer-Landau R. (2014). The moral fixed points: new di-  
855 rections for moral nonnaturalism. *Philosophical Studies* 171: 399-443.
- 856 [9] Dean, W. (2014), *Models and Computability*. *Philosophia Mathematica*  
857 22 (2): 143–166.
- 858 [10] Eklund, M. (2015). Intuitions, Conceptual Engineering, and Conceptual  
859 Fixed Points. In: Daly C. (ed.) *The Palgrave Handbook of Philosophical*  
860 *Methods*. Palgrave Macmillan: 363-385.
- 861 [11] Feferman, S. (1995). Penrose's Gödelian argument. *Psyche: An Inter-*  
862 *disciplinary Journal of Research on Consciousness*.
- 863 [12] Gödel, K. (193?), Undecidable Diophantine Propositions. In: Gödel, K.,  
864 *Collected Works, Volume III, Unpublished essays and lectures*, Feferman  
865 S., et al. (eds.), Oxford University Press 1995: 164–175.
- 866 [13] Gödel, K. (\*1951). \*1951 is Gödel's 1951 Gibbs lecture. Some basic the-  
867 orems on the foundations of mathematics and their implications, lecture  
868 manuscript. In: Gödel, K., *Collected Works, Volume III, Unpublished es-*  
869 *says and lectures*, Feferman S., et al. (eds.), Oxford University Press 1995:  
870 304—323.
- 871 [14] Halbach, V. & Horsten, L. (2005). Computational Structuralism.  
872 *Philosophia Mathematica* 13(2): 174–186.
- 873 [15] Hofstadter, D.R. (1979). *Gödel, Escher, Bach, and Eternal Golden*  
874 *Braid*. Basic Books.
- 875 [16] Krajewski, Stanislaw (2007). On Gödel's Theorem and Mechanism: In-  
876 consistency or Unsoundness is Unavoidable in any Attempt to 'Out-Gödel'  
877 the Mechanist. *Fundamenta Informaticae* 81: 173–181. Reprinted in *Top-*  
878 *ics in Logic, Philosophy and Foundations of Mathematics and Computer*  
879 *Science: In Recognition of Professor Andrzej Grzegorzczuk* (2008).

- 880 [17] Krajewski, Stanislaw (unpublished manuscript). On the Anti-Mechanist  
881 Arguments Based on Gödel's Theorem.
- 882 [18] Lucas, J.R. (1961). Minds, Machines and Gödel. *Philosophy* 36 (137):  
883 112-127.
- 884 [19] Lucas, J.R. (1968). Satan Stultified: A Rejoinder to Paul Benacerraf.  
885 *The Monist* 52: 145–158.
- 886 [20] Lucas, J.R. (1990). A paper to read to the Turing Conference at Brighton  
887 on April 6th, 1990. Unpublished manuscript <http://users.ox.ac.uk/~jrlucas/Godel/brighton.html>  
888
- 889 [21] Lucas, J.R. (1996). Minds, Machines and Gödel: A Retrospect. In: Mil-  
890 lican P. & Clark A. (eds.). *Machines and Thought*. Oxford University Press.
- 891 [22] Maddy P. (2007). *Second Philosophy. A Naturalistic Method*. Oxford  
892 University Press.
- 893 [23] Makovec, D. & Shapiro S. (eds.) (2019). *Friedrich Waismann. The Open*  
894 *Texture of Analytic Philosophy*. Springer.
- 895 [24] Nagel, E. & Newman J.R. (1958). *Gödel's Proof*. New York University  
896 Press.
- 897 [25] Nagel, E. & Newman J.R. (1961). Answer to Putnam. *Philosophy of*  
898 *Science*. 28: 209–211.
- 899 [26] Penrose, R. (1989). *The Emperor's New Mind: Concerning Computers,*  
900 *Minds and The Laws of Physics*. Oxford University Press.
- 901 [27] Penrose, R. (1994). *Shadows of the Mind: A Search for the Missing*  
902 *Science of Consciousness*. Oxford University Press.
- 903 [28] Piccinini, G. (2010/2017). *Computation in Physical Systems*. In: Zalta,  
904 E.N. (ed.) *The Stanford Encyclopedia of Philosophy*. Metaphysics Research  
905 Lab, Stanford University.
- 906 [29] Piccinini, G. (2015). *Physical Computation: A Mechanistic Account*.  
907 Oxford University Press.

- 908 [30] Plunkett D. Cappelen H. (2020). A Guided Tour Of Conceptual Engi-  
909 neering and Conceptual Ethics. In: Cappelen H., Plunkett D. & Burgess A.  
910 (eds.), *Conceptual Engineering and Conceptual Ethics*. Oxford University  
911 Press (2020): 1–26.
- 912 [31] Post, E. (1941). Absolutely Unsolvable Problems and Relatively Unde-  
913 cidable Propositions – Account of an Anticipation. In: Davis, M. (ed.).  
914 *The Undecidable* (1965). Raven Press
- 915 [32] Putnam, H.(1960). Minds and Machines. In: Hook S. (ed.). *Dimensions*  
916 *of Mind. A Symposium*. Collier-MacMillan.
- 917 [33] Putnam, H. (1980). Models and Reality. *Journal of Symbolic Logic*  
918 45(3): 464–482.
- 919 [34] Putnam, H. (1995). Review of *The Shadows of the Mind*. *Bulletin of the*  
920 *American Mathematical Society*. 32(2): 370–373.
- 921 [35] Quine, W.V.O. (1970). *Philosophy of Logic*. Harvard University Press.
- 922 [36] Quinon, P. & Zdanowski, K. (2007). Intended Model of Arithmetic.  
923 Argument from Tennenbaum’s Theorem. In: Cooper, S. Barry; Kent ,  
924 Thomas F.: Löwe, Benedikt & Sorbi, Andrea (eds.), *Computation and*  
925 *Logic in the Real World*, CiE: 313–317.
- 926 [37] Quinon, P. (2014), From Computability over Strings of Characters to  
927 Natural Numbers. In: Olszewski, A.; Brożek, B. & Urbańczyk, P. (eds.),  
928 *Church’s Thesis, Logic, Mind & Nature*. Copernicus Center Press: 310–  
929 330.
- 930 [38] Quinon, P. (2018). Taxonomy of Deviant Encodings. In: Manea F.,  
931 Miller R. & Nowotka D. (eds) *Sailing Routes in the World of Computation*.  
932 CiE 2018. *Lecture Notes in Computer Science* 10936 Springer: 338-348.
- 933 [39] Quinon, P. (2019). Can Church’s Thesis be Viewed as a Carnapian Ex-  
934 plication?, *Synthese*: Online First.
- 935 [40] Quinon, P. (2020). Implicit and explicit examples of the phenomenon of  
936 deviant encodings. *Journal of Studies in Logic, Grammar and Rhetoric*. In  
937 press.

- 938 [41] Rescrola, M. (2007), Church's thesis and the conceptual analysis of com-  
939 putability. *Notre Dame Journal of Formal Logic* 48 (2): 253–280.
- 940 [42] Shapiro, S. (1982). Acceptable Notation. *Notre Dame Journal of Formal*  
941 *Logic* 23(1): 14–20.
- 942 [43] Shapiro, S. (1998). Incompleteness, mechanism, and optimism. *Journal*  
943 *of Philosophical Logic* 4: 273–302.
- 944 [44] Shapiro, S. (2003). Mechanism, truth, and Penrose's new argument.  
945 *Journal of Philosophical Logic* 32: 19–42.
- 946 [45] Shapiro, S. (2013). Computability, Proof and Open-texture. In: Ol-  
947 szewski A., Wolenski J., Janusz R., Church's Thesis After 70 Years, Walter  
948 de Gruyter: 420-455.
- 949 [46] Turing, A. (1936). On Computable Numbers, with an Application to the  
950 Entscheidungsproblem. *Proceedings of the London Mathematical Society*  
951 42: 230–265; correction in (1937) 43: 544–546; reprinted in (Davis 1965:  
952 115–154); page numbers refer to the (1965) edition.
- 953 [47] Turing, A. (1950). Computing machinery and intelligence. *Mind*: 433–  
954 460.
- 955 [48] Turkle, S. (1978). *Psychoanalytic Politics: Jacques Lacan and Freud's*  
956 *French Revolution*. Basic Books.
- 957 [49] Turkle, S. (1984/2005). *The Second Self: The Second Self: Computers*  
958 *and the Human Spirit*. MIT Press.
- 959 [50] Turkle, S. (2011). *Alone Together: Why We Expect More from Tech-*  
960 *nology and Less from Each Other*. Basic Books.
- 961 [51] Turkle, S. (2015). *Reclaiming Conversation: The Power of Talk in a*  
962 *Digital Age*. Penguin Press.
- 963 [52] van Heuveln, B. (2000), *Emergence and consciousness: Explorations*  
964 *into the Philosophy of Mind via the Philosophy of Computation*. Ph.D.  
965 thesis. State University of New York at Binghamton.
- 966 [53] Wang, H. (1974). *From Mathematics to Philosophy*. Routledge and  
967 Kegan Paul.