

1 Article

2 Approximate and Situated Causality in Deep 3 Learning

4 Jordi Vallverdú^{1,*}

5 ¹ Philosophy Department - UAB; jordi.vallverdu@uab.cat

6 * Correspondence: jordi.vallverdu@uab.cat; Tel.: +345811618

7 **Abstract:** Causality is the most important topic in the history of Western Science, and since the
8 beginning of the statistical paradigm, its meaning has been reconceptualized many times. Causality
9 entered into the realm of multi-causal and statistical scenarios some centuries ago. Despite of
10 widespread critics, today Deep Learning and Machine Learning advances are not weakening
11 causality but are creating a new way of finding indirect factors correlations. This process makes
12 possible us to talk about approximate causality, as well as about a situated causality.

13 **Keywords:** causality; deep learning; machine learning; counterfactual; explainable AI; blended
14 cognition; mechanisms; system

15

16 1. Causalities in the 21st Century.

17 In classic Western Philosophies, causality was observed as an obvious observation of the divine
18 regularities which were ruling Nature. From a dyadic truth perspective, some events were true
19 while others were false, and those which were true followed strictly the Heaven's will. That
20 ontological perspective allowed early Greek philosophers (inspired by Mesopotamian, Egyptian and
21 Indian scientists) to define causal models of reality with causal relations deciphered from a single
22 origin, the arche (ἀρχή). Anaximander, Anaximenes, Thales, Plato or Aristotle, among others,
23 created different models about causality, all of them connected by the same idea: hazard or
24 nothingness was not possible. Despite those ideas were defended by atomists (who thought on a
25 Nature with both hazard and void), any trace of them was deleted from researches. On the other
26 hand, Eastern philosophers departed from the opposite ontological point of view: at the beginning
27 was the nothingness, and the only true reality is the continuous change of things [1]. For Buddhist
28 (using a four-valued logic), Hindu, Confucian or Taoist philosophers, causality was a reconstruction
29 of the human mind, which is also a non-permanent entity. Therefore, the notion of causality is
30 ontologically determined by situated perspectives about information values [2], which allowed and
31 fed different and fruitful heuristic approaches to reality [3], [4]. Such situated contexts of thinking
32 shape the ways by which people perform epistemic and cognitive tasks[5]–[7].

33 These ontological variations can be justified and fully understood once we assume the
34 Duhem-Quine Thesis, that is, that it is impossible to test a scientific hypothesis in isolation, because
35 an empirical test of the hypothesis requires one or more background assumptions (also called
36 auxiliary assumptions or auxiliary hypotheses). Therefore, the history of the idea of causality
37 changes coherently across the geographies and historical periods, entering during late 19th Century
38 into the realm of statistics and, later in 20th Century, in multi-causal perspectives [8]. The statistical
39 nature of contemporary causality has been involved into debates between schools, mainly Bayesians
40 and a broad range of frequentist variations. At the same time, the epistemic thresholds has been
41 changing, as the recent debate about statistical significance has shown, desacralizing the p -value.
42 The most recent and detailed academic debate on statistical significance was extremely detailed into
43 the #1 Supplement of the Volume 73, 209 of the journal *The American Statistician*, released in March,
44 20th 2019. But during the last decades of 20th Century and the beginning of the 21st Century,
45 computational tools have become the backbone of cutting scientific researches [9], [10]. After the

46 great advances produced by machine learning techniques (henceforth, ML), several authors have
47 asked themselves whether ML can contribute to the creation of causal knowledge. We will answer to
48 this question into next section.
49

50 2. Deep Learning, Counterfactuals, and Causality

51 Is in this context, where the statistical analysis rules the study of causal relationships, that we
52 find the attack to Machine Learning and Deep Learning as no suitable tools for the advance of causal
53 and scientific knowledge. The most known and debated arguments come from the eminent
54 statistician Judea Pearl [11], [12], and have been widely accepted. The main idea is that machine
55 learning do not can create causal knowledge because the skill of managing counterfactuals, and
56 following his exact words, [11] page 7: "Our general conclusion is that human-level AI cannot
57 emerge solely from model-blind learning machines; it requires the symbiotic collaboration of data
58 and models. Data science is only as much of a science as it facilitates the interpretation of data – a
59 two-body problem, connecting data to reality. Data alone are hardly a science, regardless how big
60 they get and how skillfully they are manipulated". What he is describing is the well-known problem
61 of the black box model: we use machines that process very complex amounts of data and provide
62 some extractions at the end. As has been called, it is a GIGO (Garbage In, Garbage Out) process [13],
63 [14]. It could be affirmed that GIGO problems are computational versions of Chinese room mental
64 experiment [15]: the machine can find patterns but without real and detailed causal meaning. This
65 is what Pearl means: the blind use of data for establishing statistical correlations instead that of
66 describing causal mechanisms. But is it true? In a nutshell: not at all. I'll explain the reasons.
67

68 *2.1. Deep Learning is not a data driven but a context driven technology: made by humans for humans.*

69 Most of epistemic criticisms against AI are always repeating the same idea: machines are still
70 not able to operate like humans do [16], [17]. The idea is always the same: computers are operating
71 with data using a blind semantic perspective that makes not possible that they understand the causal
72 connections between data. It is the definition of a black box model [18], [19]. But here happens the
73 first problem: deep learning (henceforth, DL) is not the result of automated machines creating by
74 themselves search algorithms and after it, evaluating them as well as their results. DL is designed by
75 humans, who select the data, evaluate the results and decide the next step into the chain of possible
76 actions. At epistemic level, is under human evaluation the decision about how to interpret the
77 validity of DL results, a complex, but still only, technique [20]. But even last trends in AGI design
78 include causal thinking, as DeepMind team has recently detailed [21], and with explainable
79 properties. The exponential growth of data and their correlations has been affecting several fields,
80 especially epidemiology [22], [23]. Initially it can be expressed by the agents of some scientific
81 community as a great challenge, in the same way that astronomical statistics modified the
82 Aristotelian-Newtonian idea of physical cause, but with time, the research field accepts new ways of
83 thinking. Consider also the revolution of computer proofs in mathematics and the debates that these
84 techniques generated among experts [24], [25].

85 In that sense, DL is just providing a situated approximation to reality using correlational
86 coherence parameters designed by the communities that use them. It is beyond the nature of any
87 kind of machine learning to solve problems only related to human epistemic envisioning: let's take
88 the long, unfinished, and even disgusting debates among the experts of different of statistical
89 schools [8]. And this is true because data do not provide or determine epistemology, in the same
90 sense that groups of data do not provide the syntax and semantics of the possible organization
91 systems to which they can be assigned. Any connection between the complex dimensions of any
92 event expresses a possible epistemic approach, which is a (necessary) working simplification. We
93 cannot understand the world using the world itself, in the same way that the best map is not a 1:1
94 scale map, as Borges wrote (1946, *On Exactitude in Science*): "...In that Empire, the Art of Cartography
95 attained such Perfection that the map of a single Province occupied the entirety of a City, and the
96 map of the Empire, the entirety of a Province. In time, those Unconscionable Maps no longer

97 satisfied, and the Cartographers Guilds struck a Map of the Empire whose size was that of the
98 Empire, and which coincided point for point with it. The following Generations, who were not so
99 fond of the Study of Cartography as their Forebears had been, saw that that vast Map was Useless,
100 and not without some Pitilessness was it, that they delivered it up to the Inclemencies of Sun and
101 Winters. In the Deserts of the West, still today, there are Tattered Ruins of that Map, inhabited by
102 Animals and Beggars; in all the Land there is no other Relic of the Disciplines of Geography”

103 Then, DL cannot follow a different information processing process, a specific one completely
104 different from those run by humans. As any other epistemic activity, DL must include different
105 levels of uncertainties if we want to use it [26]. Uncertainty is a reality for any cognitive system, and
106 consequently, DL must be prepared to deal with it. Computer vision is a clear example of that set of
107 problems [27]. Kendall and Gal have even coined new concepts to allow introduce uncertainty into
108 DL: homocedastic, and heterocedastic uncertainties (both aleatoric) [28]. The way used to integrate
109 such uncertainties can determine the epistemic model (which is a real cognitive algorithmic
110 extension of ourselves). For example, Bayesian approach provides an efficient way to avoid
111 overfitting, allow to work with multi-modal data, and make possible use them in real-time scenarios
112 (as compared to Monte Carlo approaches) [29]; or even better, some authors are envisioning
113 Bayesian Deep Learning [30]. Dimensionality is a related question that has also a computational
114 solution, as Yosuhua Bengio has been exploring during last decades [31]–[33].

115 In any case, we cannot escape from the informational formal paradoxes, which were
116 well-known at logical and mathematical level once Gödel explained them; they just emerge in this
117 computational scenario, showing that artificial learnability can also be undecidable [34]. Machine
118 learning is dealing with a rich set of statistical problems, those that even at biological level are
119 calculated at approximate levels [35]–[37], a heuristics that are being implemented also into
120 machines. This open range of possibilities, and the existence of mechanisms like informational
121 selection procedures (induction, deduction, abduction), makes possible to use DL in a controlled but
122 creative operational level [38].

123

124

125 *2.2. Deep learning is already running counterfactual approaches.*

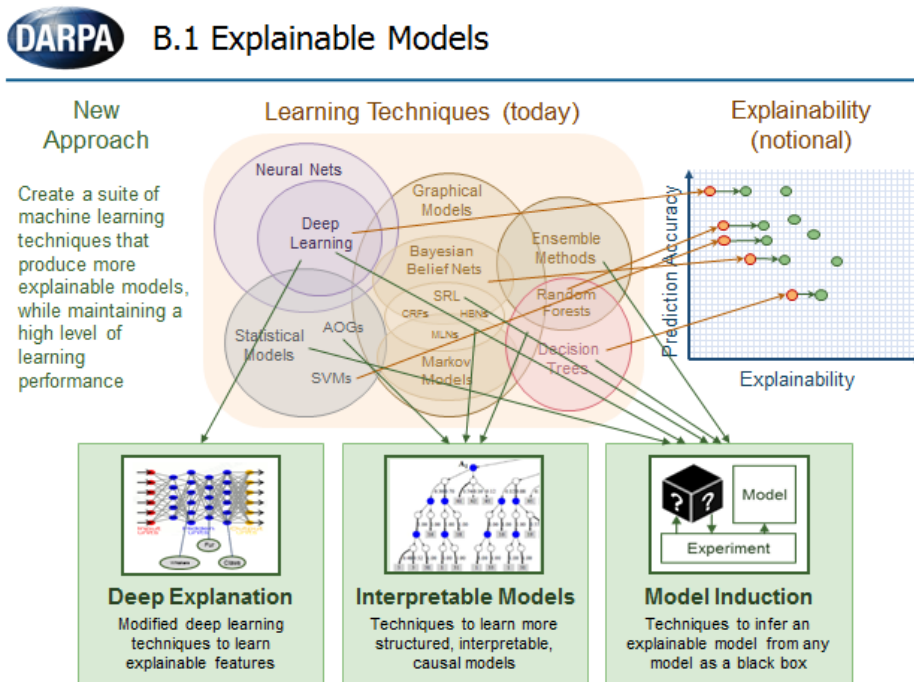
126 The second big Pearl critics is against DL because of its incapacity of integrating counterfactual
127 heuristics. First we must affirm that counterfactuals do not warrant with precision any epistemic
128 model, just add some value (or not). From a classic epistemic point of view, counterfactuals do not
129 provide a more robust scientific knowledge: a quick look at the last two thousands of both Western
130 and Eastern sciences can give support to this view [4]. Even going beyond, I affirm that counterfactual
131 can block thinking once it is structurally related to a close domain or paradigm of well-established
132 rules; otherwise is just fiction or an empty mental experiment. Counterfactuals are a fundamental
133 aspect of human reasoning [39]–[42], and their algorithmic integration is a good idea [43]. But at the
134 same time, due to the underdetermination [44]–[46], counterfactual thinking can express completely
135 wrong ideas about reality. DL can no have an objective ontology that allows it to design a perfect
136 epistemological tool: because of the huge complexity of the involved data as well as for the necessary
137 situatedness of any cognitive system. Uncertainty would not form part of such counterfactual
138 operationability [47], once it should ascribed to any not-well known but domesticable aspect of
139 reality; nonetheless, some new ideas do not fit with the whole set of known facts, the current
140 paradigm, nor the set of new ones. This would position us into a sterile no man’s land, or even block
141 any sound epistemic movement. But humans are able to deal with it, going even beyond [48].
142 Opportunistic blending, and creative innovating are part of our set of most valuable cognitive skills
143 [49].

144

145

146 *2.3. DL is not Magic Algorithmic Thinking (MAT).*

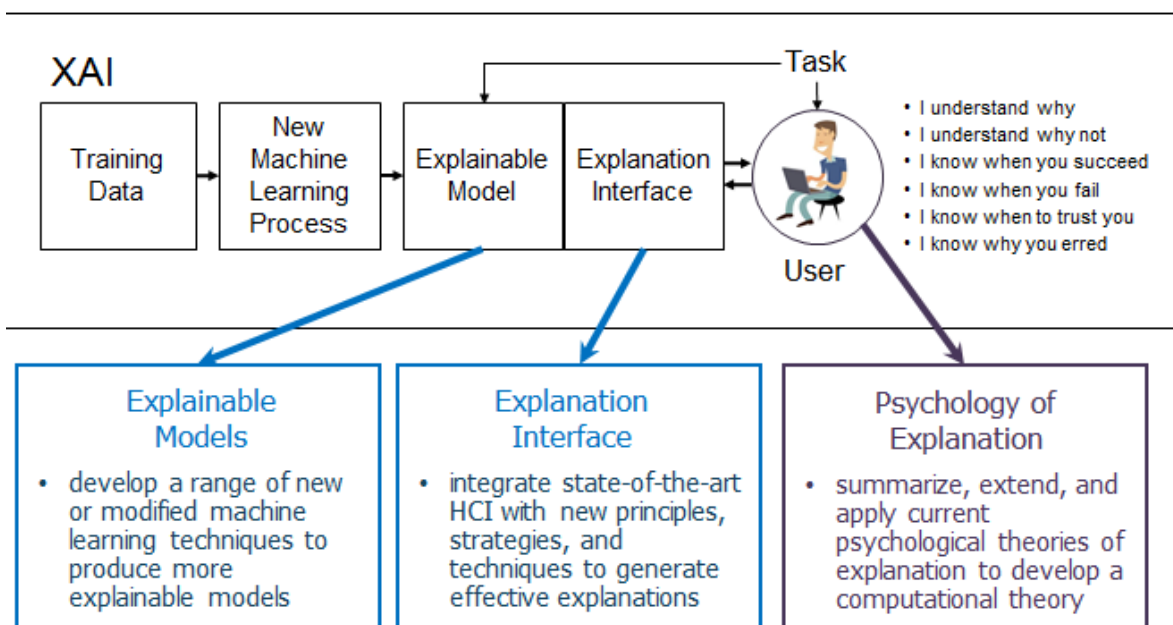
147 Our third and last analysis of DL characteristics is related to its explainability. Despite of the
 148 evidence that the causal debate is beyond any possible resolution provided by DL, because it
 149 belongs to ontological perspectives that require a different holistic analysis, it is clear that the results
 150 provided my DL must be not only coherent but also explainable, otherwise we would be in front of a
 151 new algorithmic form of magic thinking. By the same reasons by which DL cannot be just mere a
 152 complex way of curve fitting, it cannot become a fuzzy domain beyond human understanding. Some
 153 attempts are being held to prevent us from this, most of them rules by DARPA: Big Mechanisms [50]
 154 or XAI (eXplainable Artificial Intelligence) [51], [52]. An image from DARPA website:
 155



156
 157

Figure 1. Explainability in DL

DARPA B. Program Scope – XAI Development Challenges



158
 159
 160

Figure 2. XAI from DARPA

161 Again, these approaches answer to a request: how to adapt new epistemic tools to our cognitive
162 performative thresholds and characteristics. There is not a bigger revolution from a conceptual
163 perspective than the ones that happened during Renaissance with the use of telescopes or
164 microscopes. DL systems are not running by themselves interacting with the world, automatically
165 selecting the informational events to be studied, or evaluating them in relation to a whole universal
166 paradigm of semantic values. Humans neither do.

167 Singularity debates are useful and exploring possible conceptual frameworks and must be held
168 [53]–[55], but at the same time cannot become fallacious fatalist arguments against current
169 knowledge. Today, DL is a tool used by experts in order to map new connections between sets of
170 data. Epistemology is not automated process, despite of minor and naïve attempts to achieve it
171 [56], [57]. Knowledge is a complex set of explanations related to different systems that is integrated
172 dynamically by networks of epistemic (human, still) agents who are working with AI tools.
173 Machines could postulate their own models, true, but the mechanisms to verify or refine them
174 would not be beyond any mechanism different from the used previously by humans: data do not
175 express by itself some pure nature, but offers different system properties that need to be classified in
176 order to obtain knowledge. And this obtaining is somehow a creation based on the epistemic and
177 body situatedness of the system.

178
179
180
181
182
183
184
185
186
187
188
189

190 3. Extending bad and/or good human cognitive skills through DL.

191 It is beyond any doubt that DL is contributing to improve the knowledge in several areas some
192 of them very difficult to interpret because of the nature of obtained data, like neuroscience [58].
193 These advances are expanding the frontiers of verifiable knowledge beyond classic human
194 standards. But even in that sense, they are still explainable. Anyhow, are humans who fed up DL
195 systems with scientific goals, provide data (from which to learn patterns) and define quantitative
196 metrics (in order to know how close are you getting to success). At the same time, are we sure that is
197 not our biased way to deal with cognitive processes that mechanism that allows us to be creative?
198 For such reason, some attempts to reintroduce human biased reasoning into machine learning are
199 being explored [59]. This re-biasing [60], even replying emotional like reasoning mechanisms [61],
200 [62].

201 My suggestion is that after great achievements following classic formal algorithmic approaches,
202 it now time for DL practitioners to expand horizons looking into the great power of cognitive biases.

203 For example, machine learning models with human cognitive biases are already capable of
204 learning from small and biased datasets [63]. This process reminds the role of Student test in relation
205 to frequentist ideas, always requesting large sets of data until the creation of the *t*-test, but now in the
206 context of machine learning.

207 In [63] the authors developed a method to reduce the inferential gap between human beings
208 and machines by utilizing cognitive biases. They implemented a human cognitive model into
209 machine learning algorithms and compared their performance with the currently most popular
210 methods, naïve Bayes, support vector machine, neural networks, logistic regression, and random
211 forests. This even could make possible one-shot learning systems [64]. Approximate computing can

212 boost the potentiality of DL, diminishing the computational power of the systems as well as adding
213 new heuristic approaches to information analysis [35], [65]–[67].

214 Finally, a completely different type of problems, but also important, are how to reduce the
215 biased datasets or heuristics we provide to our DL systems [68] as well as how to control the biases
216 that make us not to interpret DL results properly [69]. Obviously, if there is any malicious value
217 related to such bias, it must be also controlled [70].
218

219 4. Causality in DL: the epidemiological case study

220 Several attempts has been implemented in order to allow causal models in DL, like [71] and the
221 Structural Causal Model (SCM) (as an abstraction over a specific aspect of the CNN. We also
222 formulate a method to quantitatively rank the filters of a convolution layer according to their
223 counterfactual importance), or Temporal Causal Discovery Framework (TCDF, a deep learning
224 framework that learns a causal graph structure by discovering causal relationships in observational
225 time series data) by [72]. But my attempt here will be twofold: (1) first, to consider about the value of
226 “causal data” for epistemic decisions in epidemiology; and (2) second, to look how DL could fit or
227 not with those causal claims into the epidemiological field.

228

229 4.1. Do causality affects at all epidemiological debates?

230 According to the field reference [73], MacMahon and Pugh [74] created the one of the most
231 frequently used definitions of epidemiology: “Epidemiology is the study of the distribution and
232 determinants of disease frequency in man”. Note the absence of the term ‘causality’ and, instead, the
233 use of the one of ‘determinant’. This is the result of the classic prejudices of Hill in his paper of 1965:
234 *“I have no wish, nor the skill, to embark upon philosophical discussion of the meaning of ‘causation’. The ‘cause’*
235 *of illness may be immediate and direct; it may be remote and indirect underlying the observed association. But*
236 *with the aims of occupational, and almost synonymous preventive, medicine in mind the decisive question is*
237 *where the frequency of the undesirable event B will be influenced by a change in the environmental feature A.*
238 *How such a change exerts that influence may call for a great deal of research, However, before deducing*
239 *‘causation’ and taking action we shall not invariably have to sit around awaiting the results of the research. The*
240 *whole chain may have to be unraveled or a few links may suffice. It will depend upon circumstances.”* After
241 this philosophical epistemic positioning, Hill numbered his 9 general qualitative association factors,
242 also commonly called “Hill’s criteria” or even, which is frankly sardonic, “Hill’s Criteria of
243 Causation”. For such epistemic reluctances epidemiologists abandoned the term “causation” and
244 embraced other terms like “determinant” [75], “determining conditions” [76], or “active agents of
245 change” [77]. For that reason, recent researches have claimed for a pluralistic approach to such
246 complex analysis [78]. As a consequence we can see that even in a very narrow specialized field like
247 epidemiology the meaning of cause is somehow fuzzy. Once medical evidences showed that
248 causality was not always a mono-causality [22], [79] but, instead, the result of the sum of several
249 causes/factors/determinants, the necessity of clarifying multi-causality emerged as a first-line
250 epistemic problem. It was explained as a “web of causation” [80]. Some debates about the logics of
251 causation and some popperian interpretations were held during two decades [81], [82]. Pearl himself
252 provided a graphic way to adapt human cognitive visual skills to such new epidemiological
253 multi-causal reasoning [83], as well do-calculus [84], [85], and directed acyclic graphs (DAGs) are
254 becoming a fundamental tool [86], [87]. DAGs are commonly related to randomized controlled trials
255 (RCT) for assessing causality. But RCT are not a Gold Standard beyond any critic [88], [89], because
256 as [90] affirmed, RCT are often flawed, mostly useless, although clearly indispensable (it is not so
257 uncommon that the same author claim against classic p-value suggesting a new 0,005, [91]). Krauss
258 has even defended the impossibility of using RCT without biases [92], although some authors
259 defend that DAGs can reduce RCT biases [93].

260 But there a real case that can shows us a good example about the weight of causality in real
261 scientific debates. We will see the debates about the relation between smoke and lung cancer. As
262 soon as in 1950 were explained the causal connections between smoking and lung cancer [94]. But far
263 from being accepted, these results were replied by tobacco industry using scientific experimental
264 regression. Perhaps the most famous generator of silly counterarguments was R.A. Fisher, the most
265 important frequentist researcher of 20th Century. In 1958 he published a paper in *Nature* journal [95]
266 in which he affirmed that all connections between tobacco smoking and lung cancer were due to a
267 false correlation. Even more: with the same data could be inferred that “smoking cigarettes was a
268 cause of considerable prophylactic value in preventing the disease, for the practice of inhaling is rare
269 among patients with cancer of the lung that with others” (p. 596). Two years later he was saying
270 similar silly things in a high rated academic journal [96]. He even affirmed that Hill tried to plant
271 fear into good citizens using propaganda, and entering misleadingly into the thread of
272 overconfidence. The point is: did have Fisher real epistemic reasons for not to accepting the huge
273 amount of existing causal evidences against tobacco smoking? No. And we are not affirming the
274 consequent after collecting more data not available during Fisher life. He has strong causal
275 evidences but he did not wanted to accept them. Still today, there are evidences that show how
276 causal connections are field biased, again with tobacco or the new e-cigarettes [97]–[99].

277 As a section conclusion can be affirmed that causality has strong specialized meanings and can
278 be studied under a broad range of conceptual tools. The real example of tobacco controversies offers
279 such long temporal examples.

280

281 *4.2.Can DL be of some utility for the epidemiological debates on causality?*

282 The second part of my argumentation will try to elucidate whether DL can be useful for the
283 resolution of debates about causality in epidemiological controversies. The answer is easy and clear:
284 yes. But it is directly related to a specific idea of causality as well as of a demonstration. For example,
285 can be found a machine learning approach to enable evidence based oncology practice [100]. Thus,
286 digital epidemiology is a robust update of previous epidemiological studies [101][102]. The new
287 possibilities of finding new causal patterns using bigger sets of data is surely the best advantages of
288 using DL for epidemiological purposes [103]. Besides, such data are the result of integrating
289 multimodal sources, like visual combined with classic informational [104], but the future with mode
290 and more data capture devices could integrate smell, taste, movements of agents,...deep
291 convolutional neural networks can help us, for example, to estimate environmental exposures using
292 images and other complementary data sources such as cell phone mobility and social media
293 information. Combining fields such as computer vision and natural language processing, DL can
294 provide the way to explore new interactions still opaque to us [105], [106].

295 Despite of the possible benefits, it is also true that the use of DL in epidemiological analysis has
296 a dangerous potential of unethicity, as well as formal problems [107], [108]. But again, the
297 evaluation of involved expert agents will evaluate such difficulties as things to be solved or huge
298 obstacles for the advancement of the field.

299

300

301 **5. Conclusion: causal evidence is not a result, but a process.**

302 Author has made an overall reply to main critics to Deep Learning (and machine learning) as a
303 reliable epistemic tool. Basic arguments of Judea Pearl have been criticized using real examples of
304 DL, but also making a more general epistemic and philosophical analysis. The systemic nature of
305 knowledge, also situated and even biased, has been pointed as the fundamental aspect of a new

306 algorithmic era for the advance of knowledge using DL tools. If formal systems have structural
 307 dead-ends like incompleteness, the bioinspired path to machine learning and DL becomes a reliable
 308 way [109], [110] to improve, one more time, our algorithmic approach to nature. Finally, thanks to
 309 the short case study of epidemiological debates on causality and their use of DL tools, we've seen a
 310 real implementation case of such epistemic mechanism. The advantages of DL for multi-causal
 311 analysis using multi-modal data have been explored as well as some possible critics.

312
 313

314

315

316 **Funding:** This work has been funded by the Ministry of Science, Innovation and Universities within the State
 317 Subprogram of Knowledge Generation through the research project FFI2017-85711-P Epistemic innovation: the
 318 case of cognitive sciences. This work is also part of the consolidated research network "Grup d'Estudis
 319 Humanístics de Ciència i Tecnologia" (GEHUCT) ("Humanistic Studies of Science and Technology Research
 320 Group"), recognised and funded by the Generalitat de Catalunya, reference 2017 SGR 568.

321 **Acknowledgments:** I thank Mr. Isard Boix for his support all throughout this research. Best moments are
 322 those without words, and sometimes this lack of meaningfulness entails unique meanings.

323 **Conflicts of Interest:** The author declares no conflict of interest.

324 References

325

- 326 [1] J. W. Heisig, *Philosophers of Nothingness : An Essay on the Kyoto School*. University of Hawai'i Press, 2001.
- 327 [2] J. Vallverdú, "The Situated Nature of Informational Ontologies," in *Philosophy and Methodology of*
 328 *Information*, WORLD SCIENTIFIC, 2019, pp. 353–365.
- 329 [3] M. J. Schroeder and J. Vallverdú, "Situated phenomenology and biological systems: Eastern and
 330 Western synthesis.," *Prog. Biophys. Mol. Biol.*, vol. 119, no. 3, pp. 530–7, Dec. 2015.
- 331 [4] J. Vallverdú and M. J. Schroeder, "Lessons from culturally contrasted alternative methods of inquiry
 332 and styles of comprehension for the new foundations in the study of life," *Prog. Biophys. Mol. Biol.*, 2017.
- 333 [5] A. Carstensen, J. Zhang, G. D. Heyman, G. Fu, K. Lee, and C. M. Walker, "Context shapes early
 334 diversity in abstract thought," *Proc. Natl. Acad. Sci.*, pp. 1–6, Jun. 2019.
- 335 [6] A. Norenzayan and R. E. Nisbett, "Culture and causal cognition," *Curr. Dir. Psychol. Sci.*, vol. 9, no. 4,
 336 pp. 132–135, 2000.
- 337 [7] R. E. Nisbet, *The Geography of Thought: How Asians and Westerners Think Differently...and Why: Richard E.*
 338 *Nisbett: 9780743255356: Amazon.com: Books*. New York: Free Press (Simon & Schuster, Inc.), 2003.
- 339 [8] J. Vallverdú, *Bayesians versus frequentists : a philosophical debate on statistical reasoning*. Springer, 2016.
- 340 [9] D. Casacuberta and J. Vallverdú, "E-science and the data deluge.," *Philos. Psychol.*, vol. 27, no. 1, pp.
 341 126–140, 2014.
- 342 [10] J. Vallverdú Segura, "Computational epistemology and e-science: A new way of thinking," *Minds*
 343 *Mach.*, vol. 19, no. 4, pp. 557–567, 2009.
- 344 [11] J. Pearl, "Theoretical Impediments to Machine Learning," *arXiv Prepr.*, 2018.
- 345 [12] J. Pearl and D. Mackenzie, *The book of why : the new science of cause and effect*. Basic Books, 2018.
- 346 [13] S. Hillary and S. Joshua, "Garbage in, garbage out (How purportedly great ML models can be screwed
 347 up by bad data)," in *Proceedings of Blackhat 2017*, 2017.
- 348 [14] I. Askira Gelman, "GIGO or not GIGO," *J. Data Inf. Qual.*, 2011.
- 349 [15] J. Moyal, "The Chinese room argument," in *John searle*, 2003.
- 350 [16] H. L. Dreyfus, *What Computers Can't Do: A Critique of Artificial Reason*. 1972.

- 351 [17] H. L. Dreyfus, S. E. Drey-fus, and L. a. Zadeh, "Mind over Machine: The Power of Human Intuition and
352 Expertise in the Era of the Computer," *IEEE Expert*, vol. 2, no. 2, pp. 237–264, 1987.
- 353 [18] W. N. Price, "Big data and black-box medical algorithms," *Sci. Transl. Med.*, 2018.
- 354 [19] J. Vallverdú, "Patenting logic, mathematics or logarithms? The case of computer-assisted proofs,"
355 *Recent Patents Comput. Sci.*, vol. 4, no. 1, pp. 66–70, 2011.
- 356 [20] F. Gagliardi, "The necessity of machine learning and epistemology in the development of categorization
357 theories: A case study in prototype-exemplar debate," in *Lecture Notes in Computer Science (including*
358 *subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2009.
- 359 [21] T. Everitt, R. Kumar, V. Krakovna, and S. Legg, "Modeling AGI Safety Frameworks with Causal
360 Influence Diagrams," Jun. 2019.
- 361 [22] M. Susser and E. Susser, "Choosing a future for epidemiology: II. From black box to Chinese boxes and
362 eco-epidemiology," *Am. J. Public Health*, vol. 86, no. 5, pp. 674–677, 1996.
- 363 [23] A. Morabia, "Hume, Mill, Hill, and the sui generis epidemiologic approach to causal inference.," *Am. J.*
364 *Epidemiol.*, vol. 178, no. 10, pp. 1526–32, Nov. 2013.
- 365 [24] T. C. Hales, "Historical overview of the Kepler conjecture," *Discret. Comput. Geom.*, vol. 36, no. 1, pp. 5–
366 20, 2006.
- 367 [25] N. Robertson, D. P. Sanders, P. D. Seymour, and R. Thomas, "The Four-Colour Theorem," *J. Comb.*
368 *Theory, Ser. B*, vol. 70, pp. 2–44, 1997.
- 369 [26] Yarin Gal, Y. Gal, and Yarin Gal, "Uncertainty in Deep Learning," *Phd Thesis*, 2017.
- 370 [27] A. G. Kendall, "Geometry and Uncertainty in Deep Learning for Computer Vision," 2017.
- 371 [28] A. Kendall and Y. Gal, "What Uncertainties Do We Need in Bayesian Deep Learning for Computer
372 Vision?," Mar. 2017.
- 373 [29] J. Piironen and A. Vehtari, "Comparison of Bayesian predictive methods for model selection," *Stat.*
374 *Comput.*, 2017.
- 375 [30] N. G. Polson and V. Sokolov, "Deep learning: A Bayesian perspective," *Bayesian Anal.*, 2017.
- 376 [31] Y. Bengio and Y. Lecun, "Scaling Learning Algorithms towards AI To appear in ' Large-Scale Kernel
377 Machines '," *New York*, 2007.
- 378 [32] J. P. Cunningham and B. M. Yu, "Dimensionality reduction for large-scale neural recordings," *Nat.*
379 *Neurosci.*, vol. 17, no. 11, pp. 1500–1509, 2014.
- 380 [33] S. Bengio and Y. Bengio, "Taking on the curse of dimensionality in joint distributions using neural
381 networks," *IEEE Trans. Neural Networks*, 2000.
- 382 [34] S. Ben-David, P. Hrubeš, S. Moran, A. Shpilka, and A. Yehudayoff, "Learnability can be undecidable,"
383 *Nat. Mach. Intell.*, 2019.
- 384 [35] C. B. Anagnostopoulos, Y. Ntarladimas, and S. Hadjiefthymiades, "Situational computing: An
385 innovative architecture with imprecise reasoning," *J. Syst. Softw.*, vol. 80, no. 12 SPEC. ISS., pp. 1993–
386 2014, 2007.
- 387 [36] K. Friston, "Functional integration and inference in the brain," *Progress in Neurobiology*, vol. 68, no. 2. pp.
388 113–143, 2002.
- 389 [37] S. Schirra, "Approximate decision algorithms for approximate congruence," *Inf. Process. Lett.*, vol. 43,
390 no. 1, pp. 29–34, 1992.
- 391 [38] L. Magnani, "AlphaGo, Locked Strategies, and Eco-Cognitive Openness," *Philosophies*, 2019.
- 392 [39] A. A. Baird and J. A. Fugelsang, "The emergence of consequential thought: Evidence from
393 neuroscience," *Philosophical Transactions of the Royal Society B: Biological Sciences*. 2004.

- 394 [40] T. Gomila, *Verbal Minds : Language and the Architecture of Cognition*. Elsevier Science, 2011.
- 395 [41] J. Y. Halpern, *Reasoning About Uncertainty*, vol. 2003. 2003.
- 396 [42] N. Van Hoeck, "Cognitive neuroscience of human counterfactual reasoning," *Front. Hum. Neurosci.*,
397 2015.
- 398 [43] J. Pearl, "The algorithmization of counterfactuals," *Ann. Math. Artif. Intell.*, 2011.
- 399 [44] D. Tulodziecki, "Underdetermination," in *The Routledge Handbook of Scientific Realism*, 2017.
- 400 [45] L. Laudan, "Demystifying Underdetermination," in *Philosophy of Science. The Central Issues*, 1990.
- 401 [46] D. Turner, "Local Underdetermination in Historical Science*," *Philos. Sci.*, 2005.
- 402 [47] D. Lewis, "Counterfactual Dependence and Time's Arrow," *Noûs*, 2006.
- 403 [48] M. Ramachandran, "A counterfactual analysis of causation," *Mind*, 2004.
- 404 [49] J. Vallverdú and V. C. Müller, *Blended cognition : the robotic challenge*. Springer, 2019.
- 405 [50] A. Rzhetsky, "The Big Mechanism program: Changing how science is done," in *CEUR Workshop*
406 *Proceedings*, 2016.
- 407 [51] A. Wodecki *et al.*, "Explainable Artificial Intelligence (XAI) The Need for Explainable AI," 2017.
- 408 [52] T. Ha, S. Lee, and S. Kim, "Designing Explainability of an Artificial Intelligence System," 2018.
- 409 [53] R. Kurzweil, "The Singularity Is Near: When Humans Transcend Biology," *Book*, vol. 2011. p. 652, 2005.
- 410 [54] R. V Yampolskiy, "Leakproofing the Singularity - Artificial Intelligence Confinement Problem," *J.*
411 *Conscious. Stud.*, vol. 19, no. 1–2, pp. 194–214, 2012.
- 412 [55] J. Vallverdú, "The Emotional Nature of Post-Cognitive Singularities," S. A. Victor Callaghan, James
413 Miller, Roman Yampolskiy, Ed. Springer Berlin Heidelberg, 2017, pp. 193–208.
- 414 [56] R. D. King *et al.*, "The automation of science," *Science (80-.)*, 2009.
- 415 [57] A. Sparkes *et al.*, "Towards Robot Scientists for autonomous scientific discovery," *Automated*
416 *Experimentation*. 2010.
- 417 [58] A. Iqbal, R. Khan, and T. Karayannis, "Developing a brain atlas through deep learning," *Nat. Mach.*
418 *Intell.*, vol. 1, no. 6, pp. 277–287, Jun. 2019.
- 419 [59] D. D. Bourgin, J. C. Peterson, D. Reichman, T. L. Griffiths, and S. J. Russell, "Cognitive Model Priors for
420 Predicting Human Decisions."
- 421 [60] J. Vallverdu, "Re-embodiment cognition with the same 'biases'," *Int. J. Eng. Futur. Technol.*, vol. 15, no. 1,
422 pp. 23–31, 2018.
- 423 [61] A. Leukhin, M. Talanov, J. Vallverdú, and F. Gafarov, "Bio-plausible simulation of three monoamine
424 systems to replicate emotional phenomena in a machine," *Biologically Inspired Cognitive Architectures*,
425 2018.
- 426 [62] J. Vallverdú, M. Talanov, S. Distefano, M. Mazzara, A. Tchitchigin, and I. Nurgaliev, "A cognitive
427 architecture for the implementation of emotions in computing systems," *Biol. Inspired Cogn. Archit.*, vol.
428 15, pp. 34–40, 2016.
- 429 [63] H. Taniguchi, H. Sato, and T. Shirakawa, "A machine learning model with human cognitive biases
430 capable of learning from small and biased datasets," *Sci. Rep.*, 2018.
- 431 [64] B. M. Lake, R. R. Salakhutdinov, and J. B. Tenenbaum, "One-shot learning by inverting a compositional
432 causal process," in *Advances in Neural Information Processing Systems 27 (NIPS 2013)*, 2013.
- 433 [65] Q. Xu, T. Mytkowicz, and N. S. Kim, "Approximate Computing: A Survey," *IEEE Des. Test*, vol. 33, no.
434 1, pp. 8–22, 2016.
- 435 [66] C.-Y. Chen, J. Choi, K. Gopalakrishnan, V. Srinivasan, and S. Venkataramani, "Exploiting approximate
436 computing for deep learning acceleration," in *2018 Design, Automation & Test in Europe Conference &*

- 437 *Exhibition (DATE)*, 2018, pp. 821–826.
- 438 [67] J. Choi and S. Venkataramani, "Approximate Computing Techniques for Deep Neural Networks," in
439 *Approximate Circuits*, Cham: Springer International Publishing, 2019, pp. 307–329.
- 440 [68] M. A. Gianfrancesco, S. Tamang, J. Yazdany, and G. Schmajuk, "Potential Biases in Machine Learning
441 Algorithms Using Electronic Health Record Data," *JAMA Internal Medicine*. 2018.
- 442 [69] T. Kliegr, Š. Bahník, and J. Fürnkranz, "A review of possible effects of cognitive biases on interpretation
443 of rule-based machine learning models," Apr. 2018.
- 444 [70] B. Shneiderman, "Opinion: The dangers of faulty, biased, or malicious algorithms requires independent
445 oversight," *Proc. Natl. Acad. Sci.*, 2016.
- 446 [71] T. Narendra, A. Sankaran, D. Vijaykeerthy, and S. Mani, "Explaining Deep Learning Models using
447 Causal Inference," Nov. 2018.
- 448 [72] M. Nauta, D. Bucur, C. Seifert, M. Nauta, D. Bucur, and C. Seifert, "Causal Discovery with
449 Attention-Based Convolutional Neural Networks," *Mach. Learn. Knowl. Extr.*, vol. 1, no. 1, pp. 312–340,
450 Jan. 2019.
- 451 [73] W. Ahrens and I. Pigeot, *Handbook of epidemiology: Second edition*. 2014.
- 452 [74] B. MacMahon and T. F. Pugh, *Epidemiology; principles and methods*. Little, Brown, 1970.
- 453 [75] R. Lipton and T. Ødegaard, "Causal thinking and causal language in epidemiology: it's in the details,"
454 *Epidemiol. Perspect. Innov.*, vol. 2, p. 8, 2005.
- 455 [76] P. Vineis and D. Kriebel, "Causal models in epidemiology: Past inheritance and genetic future,"
456 *Environmental Health: A Global Access Science Source*. 2006.
- 457 [77] J. S. Kaufman and C. Poole, "Looking Back on 'Causal Thinking in the Health Sciences,'" *Annu. Rev.*
458 *Public Health*, 2002.
- 459 [78] J. P. Vandenbroucke, A. Broadbent, and N. Pearce, "Causality and causal inference in epidemiology:
460 The need for a pluralistic approach," *Int. J. Epidemiol.*, vol. 45, no. 6, pp. 1776–1786, 2016.
- 461 [79] M. Susser, "Causal thinking in the health sciences concepts and strategies of epidemiology," *Causal*
462 *Think. Heal. Sci. concepts*, 1973.
- 463 [80] N. Krieger, "Epidemiology and the web of causation: Has anyone seen the spider?," *Soc. Sci. Med.*, vol.
464 39, no. 7, pp. 887–903, 1994.
- 465 [81] M. Susser, "What is a cause and how do we know one? a grammar for pragmatic epidemiology,"
466 *American Journal of Epidemiology*. 1991.
- 467 [82] C. Buck, "Popper's philosophy for epidemiologists," *Int. J. Epidemiol.*, 1975.
- 468 [83] S. Greenland, J. Pearl, and J. M. Robins, "Causal diagrams for epidemiologic research," *Epidemiology*,
469 1999.
- 470 [84] D. Gillies, "Judea Pearl Causality: Models, Reasoning, and Inference, Cambridge: Cambridge
471 University Press, 2000," *Br. J. Philos. Sci.*, 2001.
- 472 [85] R. R. Tucci, "Introduction to Judea Pearl's Do-Calculus," Apr. 2013.
- 473 [86] S. Greenland, J. Pearl, and J. M. Robins, "Causal diagrams for epidemiologic research," *Epidemiology*,
474 vol. 10, no. 1, pp. 37–48, 1999.
- 475 [87] T. J. Vanderweele and J. M. Robins, "Directed acyclic graphs, sufficient causes, and the properties of
476 conditioning on a common effect," *Am. J. Epidemiol.*, 2007.
- 477 [88] C. Boudreau, E. C. Guinan, K. R. Lakhani, and C. Riedl, "The Novelty Paradox & Bias for Normal
478 Science: Evidence from Randomized Medical Grant Proposal Evaluations," 2016.
- 479 [89] R. M. Daniel, B. L. De Stavola, and S. Vansteelandt, "Commentary: The formal approach to quantitative

- 480 causal inference in epidemiology: Misguided or misrepresented?," *International Journal of Epidemiology*,
481 vol. 45, no. 6. pp. 1817–1829, 2016.
- 482 [90] J. P. A. Ioannidis, "Randomized controlled trials: Often flawed, mostly useless, clearly indispensable: A
483 commentary on Deaton and Cartwright," *Soc. Sci. Med.*, vol. 210, pp. 53–56, Aug. 2018.
- 484 [91] J. P. A. Ioannidis, "The proposal to lower P value thresholds to .005," *JAMA - Journal of the American
485 Medical Association*. 2018.
- 486 [92] A. Krauss, "Why all randomised controlled trials produce biased results," *Ann. Med.*, vol. 50, no. 4, pp.
487 312–322, May 2018.
- 488 [93] I. Shrier and R. W. Platt, "Reducing bias through directed acyclic graphs," *BMC Medical Research
489 Methodology*. 2008.
- 490 [94] R. Doll and A. B. Hill, "Smoking and Carcinoma of the Lung," *Br. Med. J.*, vol. 2, no. 4682, pp. 739–748,
491 Sep. 1950.
- 492 [95] R. A. FISHER, "Cancer and Smoking," *Nature*, vol. 182, no. 4635, pp. 596–596, Aug. 1958.
- 493 [96] R. A. FISHER, "Lung Cancer and Cigarettes?," *Nature*, vol. 182, no. 4628, pp. 108–108, Jul. 1958.
- 494 [97] B. Wynne, "When doubt becomes a weapon," *Nature*, vol. 466, no. 7305, pp. 441–442, 2010.
- 495 [98] C. Pisinger, N. Godtfredsen, and A. M. Bender, "A conflict of interest is strongly associated with
496 tobacco industry–favourable results, indicating no harm of e-cigarettes," *Prev. Med. (Baltim.)*, vol. 119,
497 pp. 124–131, Feb. 2019.
- 498 [99] T. Grüning, A. B. Gilmore, and M. McKee, "Tobacco industry influence on science and scientists in
499 Germany.," *Am. J. Public Health*, vol. 96, no. 1, pp. 20–32, Jan. 2006.
- 500 [100] N. Ramarajan, R. A. Badwe, P. Perry, G. Srivastava, N. S. Nair, and S. Gupta, "A machine learning
501 approach to enable evidence based oncology practice: Ranking grade and applicability of RCTs to
502 individual patients.," *J. Clin. Oncol.*, vol. 34, no. 15_suppl, pp. e18165–e18165, May 2016.
- 503 [101] M. Salathé, "Digital epidemiology: what is it, and where is it going?," *Life Sci. Soc. Policy*, vol. 14, no. 1,
504 p. 1, Dec. 2018.
- 505 [102] E. Velasco, "Disease detection, epidemiology and outbreak response: the digital future of public health
506 practice," *Life Sci. Soc. Policy*, vol. 14, no. 1, p. 7, Dec. 2018.
- 507 [103] C. Bellinger, M. S. Mohamed Jabbar, O. Zaïane, and A. Osornio-Vargas, "A systematic review of data
508 mining and machine learning for air pollution epidemiology," *BMC Public Health*. 2017.
- 509 [104] S. Weichenthal, M. Hatzopoulou, and M. Brauer, "A picture tells a thousand...exposures:
510 Opportunities and challenges of deep learning image analyses in exposure science and environmental
511 epidemiology," *Environment International*. 2019.
- 512 [105] G. Eraslan, Ž. Avsec, J. Gagneur, and F. J. Theis, "Deep learning: new computational modelling
513 techniques for genomics," *Nat. Rev. Genet.*, vol. 20, no. 7, pp. 389–403, Jul. 2019.
- 514 [106] M. J. Cardoso *et al.*, Eds., *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical
515 Decision Support*, vol. 10553. Cham: Springer International Publishing, 2017.
- 516 [107] C. Kreatsoulas and S. V Subramanian, "Machine learning in social epidemiology: Learning from
517 experience.," *SSM - Popul. Heal.*, vol. 4, pp. 347–349, Apr. 2018.
- 518 [108] "In the age of machine learning randomized controlled trials are unethical." [Online]. Available:
519 [https://towardsdatascience.com/in-the-age-of-machine-learning-randomized-controlled-trials-are-unet
520 hical-74acc05724af](https://towardsdatascience.com/in-the-age-of-machine-learning-randomized-controlled-trials-are-unethical-74acc05724af). [Accessed: 03-Jul-2019].
- 521 [109] T. F. Drumond, T. Viéville, and F. Alexandre, "Bio-inspired Analysis of Deep Learning on Not-So-Big
522 Data Using Data-Prototypes," *Front. Comput. Neurosci.*, vol. 12, p. 100, Jan. 2019.

- 523 [110] K. Charalampous and A. Gasteratos, "Bio-inspired deep learning model for object recognition," in 2013
524 *IEEE International Conference on Imaging Systems and Techniques (IST)*, 2013, pp. 51–55.
525