

dialectica

International Journal of Philosophy

Contents

LÁSZLÓ BERNÁTH & TAMÁS PAÁR, <i>Responsibility First: How to Resist Agnosticism About Moral Responsibility</i>	525
BENJAMIN WINOKUR, <i>Self-Knowledge and Interpersonal Reasoning</i>	547
PAUL O. IRIKFE, <i>Avner Baz's Ordinary Language Challenge to the Philosophical Method of Cases</i>	571
DAN MARSHALL, <i>The Dual-Detector Argument Against the Modal Theory</i> . .	597
BOAZ FARADAY SCHUMAN, <i>Lewisian Worlds and Buridanian Possibilia</i>	623
E. E. SHENG, <i>Weakly Discerning Vertices in a Plenitude of Graphs</i>	651

dialectica

International Journal of Philosophy

Official Organ of the European Society of Analytic Philosophy

founded in 1947 by Gaston Bachelard, Paul Bernays and Ferdinand Gonseth

Editorial Board

Jérôme Dokic, EHESS, Paris, France

Pascal Engel, EHESS, Paris, France

Manuel García-Carpintero, Universitat de Barcelona, Spain

Diego Marconi, Università di Torino, Italy

Carlos Moya, Universitat de València, Spain

Martine Nida-Rümelin, Université de Fribourg, Switzerland

François Recanati, Collège de France, Paris

Marco Santambrogio, Università degli Studi di Parma, Italy

Peter Simons, Trinity College Dublin, Ireland

Gianfranco Soldati, Université de Fribourg, Switzerland

Marcel Weber, Université de Genève, Switzerland

Editors

Fabrice Correia, University of Geneva

Philipp Blum, University of Lucerne

Sharon Casu, University of Fribourg (managing editor)

Marco R. Schori, University of Bern (managing editor)

Review Editors

Stephan Leuenberger and Philipp Blum

Editorial Committee

Philipp Blum (né Keller), Claudio Calosi, Sharon Casu, Zoé Christoff, Fabrice Correia, Matthias Egg, Martin Glazier, Aleks Knoks, Arturs Logins, Jörg Löschke, Giovanni Merlo, Robert Michels, Ryan Miller, Jonathan Mitchell, Paolo Natali, Donnchadh O'Conaill, François Pellet, Edgar Phillips, Jan Plate, Stephanie Rennick, Thomas Schindler, Sebastian Schmidt, Marco R. Schori, David Schroeren, Mike Stuart, Daniel Vanello.

Consulting Board

Johannes Brandl (Salzburg), João Branquinho (Lisboa), Elke Brendel (Bonn), Ingar Brinck (Lunds), Eros Corazza (Ikerbasque and Carleton), Josep Corbi (València), Michael Esfeld (Lausanne), Dagfinn Føllesdal (Stanford and Oslo), Frank Jackson (Australian National University, Canberra), Max Kistler (Paris I), Max Kölbel (Wien), Jan Lacki (Genève), Karel Lambert (Irvine), Paolo Leonardi (Bologna), Fraser Macbride (Manchester), Josep Macià (Barcelona), Genoveva Martí (Barcelona), Élisabeth Pacherie (Institut Jean Nicod, Paris), David Piñeda (Girona), Wlodek Rabinowicz (Lund), Barry Smith (Buffalo), Christine Tappolet (Montréal), Neil Tennant (Ohio State), Mark Textor (King's College London), Achille Varzi (Columbia University), Alberto Voltolini (Torino), Timothy Williamson (Oxford).

December 2022

Contents

LÁSZLÓ BERNÁTH & TAMÁS PAÁR, <i>Responsibility First: How to Resist Agnosticism About Moral Responsibility</i>	525
BENJAMIN WINOKUR, <i>Self-Knowledge and Interpersonal Reasoning</i>	547
PAUL O. IRIKEFE, <i>Avner Baz's Ordinary Language Challenge to the Philosophical Method of Cases</i>	571
DAN MARSHALL, <i>The Dual-Detector Argument Against the Modal Theory</i> ..	597
BOAZ FARADAY SCHUMAN, <i>Lewisian Worlds and Buridanian Possibilia</i>	623
E. E. SHENG, <i>Weakly Discerning Vertices in a Plenitude of Graphs</i>	651

PROOF

Responsibility First

How to Resist Agnosticism About Moral Responsibility

LÁSZLÓ BERNÁTH & TAMÁS PAÁR

We argue against the view that one should suspend belief in the existence of moral responsibility. We start with a simple argument based on the claim that the existence of obligations entails the existence of moral responsibility. If this is true, then agnosticism about moral responsibility is incoherent. However, this simple argument is insufficient. It can be repaired by focusing on agents who rationally believe in a particular conception of obligation (the “Responsibility First View”). On that conception, non-moral obligations that are not appropriately related to moral obligations can be freely ignored and the property of being morally responsible is identical to the property of fulfilling all necessary conditions for bearing moral obligations. Those agents who rationally hold RFV can still rationally believe in moral responsibility even if they lack direct evidence for the existence of moral responsibility.

Even if we lack evidence for the existence of moral responsibility or if scientific research makes it unlikely that moral responsibility is real, one can still rationally maintain belief in it as long as one adopts a specific view of moral obligations (the “Responsibility First View”). Or so we will argue.

We first outline the case for agnosticism about responsibility (section 1), then we sketch a simple objection against it (section 2) and the reasons why the simple objection fails (section 3). Next, we outline the Responsibility First View (section 4) and we reformulate the case against agnosticism in a form that is not subject to the earlier difficulties (section 5). Finally, we consider free will, fairness and the circle of responsible agents.

2.1 Agnosticism About Moral Responsibility

26 There are three basic epistemological stances about moral responsibility. Some
 27 believe that normal human adults are often morally responsible, others deny
 28 that we know they are. The second approach, in turn, has two distinct versions.
 29 Members of the first group maintain that nobody is ever morally responsible
 30 (see Strawson, G. 1994; Pereboom 2001; Levy 2011), and they imply that this
 31 belief is justified for all of us. Members of the second group argue that we
 32 don't have enough evidence to tell. Our evidence is not decisive with regard
 33 to the existence of moral responsibility. Call philosophers who belong to this
 34 group "agnostics about moral responsibility." Typically, they suspend belief in
 35 the existence of moral responsibility, and they think that others should join
 36 them in doing so. A number of philosophers have put forward arguments to
 37 (roughly) that effect. They do not explicitly deny the reality of responsibil-
 38 ity, but they argue that our most popular (compatibilist or incompatibilist)
 39 theories of moral responsibility make it unlikely that we could tell whether
 40 anyone is morally responsible (Byrd 2010, 2021; Sehon 2013, 2016; Kearns
 41 2015).¹

1 Note that although the papers cited here are centered around arguments that seem to support responsibility agnosticism, there are important differences between the accounts of the authors. Byrd (2010) embraces agnosticism about moral responsibility, claiming that the current debate should lead us to accept that we do not *now* know whether anyone is ever morally responsible. Nevertheless, in Byrd (2021) he seems to mitigate the strength of his argument, arguing in his conclusion that even if at present agnosticism about free will is reasonable, there seems to be hope that we can overcome our ignorance *in the future*. Kearns argues mostly for agnosticism about *free will*, but he also believes that "the thesis that we (don't) have free will [...] entails the moral claim that we are (not) morally responsible" (2015, 249), hence what he says amounts to a version of moral responsibility agnosticism too. What he insists is that we do not *know* whether there is free will or not, not that we are *unjustified* in believing it, nevertheless, by "not know" he means that our justification is so weak that it doesn't even meet a "low standard" (2015, 236). But having such a weak justification in a given question might very well warrant suspending belief. Sehon (2013, 2016) seems to be the furthest from the position of these agnostics, as he develops a certain variant of compatibilism in order to counter his own challenge against belief in moral responsibility. However, there is a good reason to consider and also to answer the agnostic *arguments* of these three authors together. The reason is that their agnostic arguments are logically independent from answers that they might come up with answering or at least mitigating moral responsibility agnosticism itself, as Sehon (2016) does. Naturally, one could be consistent in accepting their arguments that support agnosticism about moral responsibility while rejecting ways they may propose to evade this kind of agnosticism. (This is why it is no surprise that Sehon 2013 basically employs agnostic arguments, without offering a solution.) Therefore, one can scrutinize these arguments independently of the full-blown theories of the aforementioned authors. It is worth noting that their ways to avoid the agnostic conclusion seem

Those who try to resist agnosticism seek to show that we do have sufficient evidence—be it moral (van Inwagen 1983, 206–223; Coffman 2016), phenomenological (Guillon 2014), conceptual (Latham 2019), transcendental (Lockie 2018), or practice-based (Strawson, P.F. 1969)—to make our belief in the reality of moral responsibility justified. Agnostics, however, retort that such pieces of evidence are unreliable and open to objections. In this paper, we do not engage with that debate.² Rather, for the sake of argument, we take it for granted that there is no sufficient evidence for the reality of moral responsibility, and we try to show that one can rationally attribute moral responsibility to herself and others even in that case. Those who prefer to argue against responsibility agnosticism more directly, and believe that there is sufficient evidence in favor of moral responsibility, may still welcome our argument as an additional way to counter the agnostic.

The argument for agnosticism can be formulated in the following way:

(AR1) Nobody is justified in believing that the metaphysical conditions of moral responsibility are ever satisfied.

(AR2) If you are not justified in believing that a necessary condition of X 's existence is satisfied, then you are not justified in believing that X exists.

Therefore,

(AR3) Belief in the existence of moral responsibility is unjustified.

On a common interpretation of justification, the following principle is true:

to be controversial. For example, Sehon (2016) offers a non-standard, non-causal and at present unpopular blend of compatibilism that has met with serious criticism (Mele 2019). If this kind of criticism is correct, the argument of this paper may still be sound, as it offers a different way to counter agnosticism about moral responsibility. Furthermore, even if the ways to respond to moral responsibility agnosticism that are suggested by those who themselves employ agnostic arguments might work, our argument does not lose its significance: it is an additional way to counter the kind of agnosticism in question that could strengthen belief in moral responsibility even more.

² It may be worthwhile to point out that our argument is somewhat akin to the transcendental arguments for free will and responsibility such as Robert Lockie's recent arguments (2018). For example, like Lockie's argument, we argue for the rationality of believing in moral responsibility is based on some analysis of conditions for bearing obligations. However, transcendental arguments aim to show that all rational (human) agents should believe in free will and responsibility (and Lockie's transcendental arguments share the same ambitions) whereas the argument we present attempts to show only that some agents can rationally believe in moral responsibility.

63 S. We should suspend those of our beliefs that are not justified.³

64 And so, one can conclude that

65 ARC. We should suspend belief about the existence of moral re-
66 sponsibility.

67 The agnostic ascribes an epistemic obligation to those who assess the evidence
68 regarding the existence of moral responsibility. The core intuition of our paper
69 is that there is a serious tension between suspending beliefs about moral
70 responsibility and ascribing epistemic obligations to oneself and others—
71 intuitively, someone who is not morally responsible cannot have obligations.
72 One can argue for this in two ways. First, one could argue that no one can
73 be obliged to suspend belief about moral responsibility. Alternatively, one
74 can say that holding a specific conception of moral responsibility makes it
75 irrational to believe in obligations to suspend belief about moral responsibility.
76 In the next section, we explore the first idea in order to see if a simple and
77 intuitive argument could support it.

72 The Simple Objection

79 It might seem *prima facie* plausible that moral responsibility is a precondition
80 of having obligations. If it is, then ascribing obligations while suspending
81 belief about moral responsibility is irrational.

82 Why is it plausible that moral responsibility is a precondition of having
83 obligations? One could appeal to the idea that only morally responsible agents
84 can have normatively binding obligations. Consider the following argument:

85 SIMPLE OBJECTION.

86 (SO1) If nobody is morally responsible, then nobody has normatively binding
87 obligations.

88 (SO2) If nobody has normatively binding obligations, then nobody has a nor-
89 matively binding obligation to suspend any of her beliefs.

3 For example, one of the most prominent moral skeptics writes: “To call a belief ‘justified’ is to say that the believer ought to hold that belief as opposed to suspending belief, because the believer has adequate epistemic grounds for believing that it is true (at least in some minimal sense)” (Sinnott-Armstrong 2008, 48).

90 (SO3) If nobody has a normatively binding obligation to suspend any of her
 91 beliefs, then nobody has a normatively binding obligation to suspend
 92 belief in moral responsibility.

93 (SO4) If nobody is morally responsible, then nobody has a normatively binding
 94 obligation to suspend belief in moral responsibility.

95 Either there are morally responsible agents or not. By (SO4), if there are no
 96 such agents, then nobody is obliged to suspend belief in moral responsibility.
 97 On the other hand, if there *are* morally responsible agents, some of whom
 98 are obliged to suspend some of their beliefs which could not be the case were
 99 they not morally responsible, then, one might argue, nobody can have a good
 100 reason to suspend belief in moral responsibility. For—assuming that moral
 101 responsibility is a precondition of having normatively binding obligations—
 102 an agent cannot ascribe to herself an obligation to suspend belief in moral
 103 responsibility unless she also takes herself to be morally responsible. That
 104 sounds incoherent, so no one can consistently believe, in the light of (SO4),
 105 that there is an obligation to be agnostic about moral responsibility. In short:

106 RATIONALITY PREMISE. If (SO4) is true, then nobody has a norma-
 107 tively binding obligation to suspend belief in moral responsibility.

108 If the argument so far is sound, it follows that

109 (SOC) No one has a normatively binding obligation to suspend belief in moral
 110 responsibility.

111 The conclusion of the Simple Argument is a threat to agnosticism because
 112 it implies that nobody has a normatively binding obligation to suspend her
 113 belief in moral responsibility, even if there is no direct evidence for the re-
 114 ality of moral responsibility. On the other hand, if nobody is in fact morally
 115 responsible, then the lack of evidence for the existence of moral responsibility
 116 does not matter, since the lack of evidence fails to have normatively binding
 117 consequences. The agnostic's claim that we ought to suspend belief in moral
 118 responsibility is thus refuted.

1193 Why the Simple Objection Fails

120 One can challenge the **SIMPLE OBJECTION** on a number of grounds. Here we
 121 take into account two objections to (SO₁) and one against the **RATIONALITY**
 122 **PREMISE**.

123 The first problem about (SO₁) is the following. Even if being morally re-
 124 sponsible is a precondition of having moral obligations, some normatively
 125 binding obligations might not be moral in nature, and so having them does
 126 not entail being morally responsible. Indeed, normatively binding obligations
 127 come in many varieties. One might have epistemic, aesthetic, prudential, legal,
 128 as well as role obligations. Moreover, on many theories of epistemic obligation,
 129 epistemic obligations are not moral at all (see, for instance, **Feldman 1988**;
 130 **Russell 2001**). If the obligation to suspend judgment about moral responsibil-
 131 ity is a non-moral, epistemic obligation, then one could have it even if one is
 132 not morally responsible. Hence (SO₁) seems to be false.

133 Another important objection to (SO₁) is that moral responsibility may not
 134 be a precondition of having moral obligations, or so the agnostic could argue.
 135 She could rightly claim that if we conceive of moral obligations in a certain
 136 way, then it is logically possible for agents who are not morally responsible
 137 to be nonetheless morally obliged to do something. For example, one might
 138 conceive of moral obligation in a consequentialist fashion and say that we
 139 have a moral obligation to maximize pleasure and minimize suffering. And
 140 it is possible that whether or not anyone is morally responsible, suspending
 141 belief in moral responsibility would minimize the amount of suffering. Thus,
 142 it could be the case that someone bears a (consequentialist) moral obligation
 143 to suspend belief in moral responsibility regardless of whether she is a morally
 144 responsible agent (see **Smilansky 1994**; **Pereboom 2001**; **Waller 2004**; **Trakakis**
 145 **2007**).

146 Further, the **RATIONALITY PREMISE** is open to the objection that there
 147 is a gap between the truth of a proposition and rationally believing that
 148 proposition. Even if (SO₁), (SO₂), and (SO₃) are true, it does not necessarily
 149 follow that everyone is rational in believing any of those premises. An agent's
 150 epistemic position might be such that her evidence either contradicts one
 151 of (SO₁)–(SO₃) or does not justify any of them. The agent's evidence may
 152 even be such that it is rational for her to believe in the soundness of the
 153 agnostic's argument. So even if (SO₄) is true, there could be a normatively
 154 binding obligation to suspend belief in moral responsibility. Therefore, the
 155 **RATIONALITY PREMISE** appears to be false.

156 In order to avoid these difficulties, we need to modify the **SIMPLE OBJEC-**
 157 **TION**. Instead of talking about normatively binding obligations, we will focus
 158 on more specific ones. To evade the second challenge, we will base the argu-
 159 ment on a particular conception of moral responsibility, one that, if rationally
 160 upheld, renders it irrational to ascribe to oneself an obligation to be agnostic
 161 about moral responsibility. Finally, to avoid the difficulties with the **RATIO-**
 162 **NALITY PREMISE**, we will defend only those agents' beliefs who rationally
 163 accepted such a conception. The next section describes the conception that
 164 we will work with, the Responsibility First View, in detail.

164 **The Responsibility First View**

166 Consider the following famous passage from Wittgenstein:

167 Supposing that I could play tennis and one of you saw me playing
 168 and said “Well, you play pretty badly” and suppose I answered “I
 169 know, I’m playing pretty badly but I don’t want to play any better,”
 170 all the other man could say would be “Ah, then that’s all right.”
 171 But suppose I had told one of you a preposterous lie and he came
 172 up to me and said, “You’re behaving like a beast” and then I were
 173 to say “I know I behave badly, but then I don’t want to behave any
 174 better,” could he then say “Ah, then that’s all right”? Certainly not;
 175 he would say “Well, you ought to want to behave better.” (1965, 5)

176 When someone says “Well, you play pretty badly,” in most cases she is not
 177 merely offering a description but implies roughly the following: “You should
 178 do something about it if you don’t want to look ridiculous.” In Wittgenstein’s
 179 story, the player in effect replies that he does not care about this implied
 180 “should.” Using contemporary terms, one could say that the implied “should”
 181 expressed a prudential obligation to prevent an undesirable outcome, such as
 182 being ridiculed. Note that even if we suppose that the player would be unhappy
 183 if someone actually ridiculed him, he could nonetheless reply “That’s all right,”
 184 because being imprudent is not an unacceptable normative error. It seems to
 185 be implausible to think that one should avoid prudential errors with all her
 186 strength in every situation or that one should feel remorse if she made such
 187 an error. Sometimes, it is all right not to care about prudential obligations
 188 even if they actually bind the agent. In other words, it might be OK to neglect
 189 them even if violating them constitutes a basis for some negative treatment
 190 (such as ridicule).

191 However, the second example suggests that violating an obligation *is* a norma-
 192 tive error that is unacceptable to such a great extent that one should feel
 193 remorse and should avoid repeating the error with all her strength. These vio-
 194 lations are just not “all right”; they cannot be shrugged off. Wittgenstein and
 195 many other philosophers claim that moral obligations fall into this category.
 196 Violating them results in unacceptable normative errors. Further, Wittgen-
 197 stein’s paper seems to imply that *only* the violation of moral obligations results
 198 in such an error. We will call this idea the Moral Primacy Thesis (MPT).

199 MPT is central to our case, so we would like to express it more precisely
 200 (incidentally explaining why the term “moral primacy” is apt). The following
 201 definition of “all things considered obligations” will be useful for that purpose:

202 ALL THINGS CONSIDERED OBLIGATION TO $\varphi =_{df}$. An obligation
 203 which is not overridden by any other obligation (in the given case)
 204 and which prescribes doing φ to agent *S* in a way that *S* should avoid
 205 violating the obligation with all her strength; and if *S* fails to observe
 206 the obligation to φ , then *S* should feel remorse.

207 We follow here Searle (1978) and many other philosophers who used the
 208 term “all things considered obligation.” Nevertheless, we add that a genuine
 209 all things considered obligation to perform a specific action must have a
 210 normative weight that makes its violation normatively unacceptable. If an
 211 obligation does not have the significant normative weight, then—all things
 212 considered—it is permissible to ignore it.

213 Philosophers often talk about obligations that have a tendency to constitute
 214 all things considered obligations. They call these *pro tanto* or *prima facie*
 215 obligations (Ross 1930). These tend to constitute all things considered obliga-
 216 tions if other, stronger obligations do not override them. (The paradigmatic
 217 examples are moral obligations.) However, for our purposes, it is better to not
 218 commit ourselves to any specific understanding of *pro tanto* or *prima facie*
 219 obligations because not only the difference between *pro tanto* and all things
 220 considered obligations is relevant for our argument but the difference between
 221 obligations that can constitute all things considered obligations in themselves
 222 and obligations that can do this only by the help of other obligations. So,
 223 instead of talking about *prima facie* and *pro tanto* obligations, we will use the
 224 term “strong obligation,” defined as follows:

225 STRONG OBLIGATION TO $\varphi =_{df}$. An obligation that constitutes an
 226 all things considered obligation to φ (in the given case) unless it is
 227 overridden by some other strong obligation(s) to do something else.

228 So in some cases, strong obligations to φ constitute an all things considered
 229 obligation to φ , and in other cases, strong obligations to φ do not constitute
 230 an all things considered obligation to φ (if they are overridden by other strong
 231 obligations).⁴

232 In addition, there are obligations that fail to constitute all things considered
 233 obligations in spite of the fact that nothing overrides them. For instance, in
 234 many cases, prudential obligations do not constitute all things considered
 235 obligations even though the agent has no other kind of obligation. This is pre-
 236 cisely the case in Wittgenstein's example: although the agent has a prudential
 237 obligation to play tennis better, he is free to ignore and violate it. We call these
 238 obligations *weak* obligations.

239 WEAK OBLIGATION TO $\varphi =_{df}$. An obligation that does not constitute
 240 an all things considered obligation to φ unless it appropriately relates
 241 to a strong obligation in the given case.

242 We intentionally use the vague term "appropriately relates." It is a complicated
 243 question when and how strong obligations turn weak obligations into all
 244 things considered obligations. For the present purposes, what matters is that
 245 this is certainly possible—whatever the details. For instance, if the tennis
 246 player in Wittgenstein's example had previously promised his wife to do his
 247 best and avoid ridicule, and there was no strong obligation to override the
 248 obligation to keep his promise, then he would have an all things considered
 249 obligation to play better. In this case, his prudential obligation to play better
 250 would be an all things considered obligation, because it would be appropriately
 251 related to his moral obligation to fulfill his promise to avoid ridicule.

252 Using the terminology just introduced, we can now characterize MPT more
 253 precisely:

4 It is worthwhile to note that our notion of strong obligation resembles the Kantian notion of categorical imperative. The main difference is that our notion of strong obligation does not imply universalizability. That is, we do not deny the possibility that one may have an obligation that constitutes an all things considered obligation to φ , although one cannot at the same time will that it becomes a universal law.

254 MORAL PRIMACY THESIS (MPT). All moral obligations are strong
 255 obligations and every other kind of obligation is weak.

256 In other words, only moral obligations can constitute all things considered
 257 obligations without being appropriately related to other kinds of obligations.
 258 On the other hand, prudential, epistemic, role, legal, etc. obligations can only
 259 constitute all things considered obligations if they are appropriately related
 260 to moral obligations.

261 It is easy to see that one of the relevant consequences of **MPT** is the follow-
 262 ing:

263 **NORMATIVE WEAKNESS OF NON-MORAL OBLIGATIONS.** Nobody
 264 has to avoid violating with all her strength those non-moral obliga-
 265 tions that do not relate appropriately to any of her moral obligations
 266 and if someone fails to observe such an obligation, she should not
 267 feel remorse.

268 Attributing moral obligation to agents must have some conditions. For exam-
 269 ple, it would certainly be absurd to attribute moral obligations to beings that
 270 are incapable to act, because it makes no sense to say that they should avoid
 271 doing something with all their strength. Whatever the relevant conditions
 272 are, there is an obvious term, under **MPT**, for those beings who fulfill all of
 273 them: they are the morally responsible agents. So proponents of **MPT** are free
 274 to adopt the following thesis as a component of their moral framework:

275 **RESPONSIBILITY IDENTITY THESIS (RIT).** The property of being
 276 morally responsible is identical to the property of fulfilling all con-
 277 ditions for bearing moral obligations.

278 **RIT** makes moral responsibility into a precondition of having moral obliga-
 279 tions, and it also makes the former prior to the latter in a certain respect. Since
 280 moral obligations are, in turn, prior to any other type of obligations under
 281 **MPT**, we will call the combination of **MPT** and **RIT** the Responsibility First
 282 View (RFV).

285 5 The Primacy Argument

284 We can now turn to the revision of the **SIMPLE OBJECTION**. The proponent of
 285 **RFV** can answer the agnostic's challenge as follows: "You claim that I should
 286 suspend my belief in moral responsibility because I cannot prove that anyone
 287 meets the conditions for being morally responsible. However, based on my
 288 conception of morality and responsibility, if I am not morally responsible, then
 289 I do not have any moral obligations. And if I have no moral obligations, I do not
 290 have any obligations that I should fulfill with all my strength, any obligations
 291 that should seriously concern me. In technical terms, I do not have all things
 292 considered obligations. So if I am not morally responsible, then it is all right
 293 for me to disregard your demand about suspension of belief. And in case you
 294 claimed that I have an all things considered obligation to suspend my belief,
 295 an obligation which I cannot disregard without committing a normative fault
 296 I should regret, then I conclude on the basis of my conception of responsibility
 297 that I am a morally responsible being after all. Either way, I can rationally
 298 resist your challenge and keep believing in moral responsibility."

299 We would like to express this revised version of the **SIMPLE OBJECTION**
 300 more formally:

301 PRIMACY ARGUMENT.

- 302 (PA1) No agent can rationally think that she has an all things considered
 303 obligation to suspend her belief that she fulfills the necessary conditions
 304 of having all things considered obligations.
- 305 (PA2) Someone who rationally upholds **RFV** cannot rationally think that she
 306 has an all things considered obligation to suspend her belief that she is
 307 morally responsible.
- 308 (PA3) If someone cannot rationally think that she has an all things considered
 309 obligation to suspend her belief that she is morally responsible, then
 310 she is rational to reject agnosticism about moral responsibility.
- 311 (PAC) Rejecting agnosticism about moral responsibility is rational for anyone
 312 who rationally upholds **RFV**.

313 Until the agnostic does not challenge the moral framework that proponents of
 314 **RFV** employ, she cannot undermine their belief in their own moral responsi-
 315 bility. The agnostic cannot challenge belief in moral responsibility by merely
 316 pointing out that evidence for the existence of moral responsibility is insuf-
 317 ficient. What is more, if someone upholds **RFV** rationally, then it would be

318 straightforward irrational for her to accept the agnostic's conclusion—unless
319 she finds out that her own moral framework is untenable.

320 The proponent of **RFV** gains a huge dialectical advantage by deploying the
321 **PRIMACY ARGUMENT**. Due to the **PRIMACY ARGUMENT**, the debate shifts
322 from the sufficiency of evidence to the tenability of a specific moral framework.
323 Defending the tenability of **RFV** seems to be much easier than defending
324 the sufficiency of evidence regarding the existence of moral responsibility.
325 Especially so if the proponent of the **MORAL PRIMACY THESIS**, by investigating
326 the nature of moral obligation and responsibility, comes to the conclusion
327 that moral responsibility has heavy-weight metaphysical preconditions such
328 as libertarian free will, since scientific evidence for libertarian freedom seems
329 to be lacking. (We will say more about free will in the last section.)

330 Moreover, as far as we can tell, both **MPT** and **RIT** can be supported by
331 considerable arguments. Even though some consequentialists deny that moral
332 obligation implies moral responsibility, that principle seems to be fundamental
333 and obvious for almost everyone—as even consequentialist critics note (**Waller**
334 **2004, 427–428**). And someone who upholds **RIT** can explain why that principle
335 is true: being a morally responsible agent is the same as being a potential
336 bearer of moral obligations.

337 **MPT** also has notable advantages. Many people find it plausible that moral
338 obligations can override all other obligations. **MPT** explains why this is the
339 case: the set of moral obligations is identical to the set of strong obligations.
340 Additionally, **MPT** provides a substantive definition of moral obligation: moral
341 obligations are those obligations that can constitute all things considered
342 obligation without the involvement of other types of obligation. Another
343 notable advantage of **MPT** is that it helps understanding why some obligations
344 can be neglected without normative costs in certain cases but not in others.
345 Take, for example, the highway code, which prescribes various patterns of
346 behaviour (call them “legal obligations”). Some of those prescriptions can be
347 non-culpably neglected in a completely abandoned city. Still, in most cases,
348 violating them is normatively unacceptable. One can use **MPT** to explain this
349 phenomenon by pointing out that the highway code contains weak obligations.
350 In most cases, they are appropriately related to moral obligations (for example,
351 to the obligation to secure the safety of human beings). However, they are not
352 appropriately related to moral obligations in an abandoned city.

353 Of course, anyone, including the agnostic, can argue against **RFV**. Indeed,
354 it seems that one can find not only prominent supporters of **RFV** (Kant seems
355 to be the most obvious example) but able critics too. For instance, Bernard

356 Williams criticizes an ethical system under the label “morality” that contains,
357 among other things, *RFV* (see Williams 2006, 174–196), because he believes
358 that moral systems with such a strong notion of moral obligations threaten
359 personal integrity. Even though the investigation of such counterarguments
360 that are based on such wide-ranging considerations about the relation between
361 a whole system of morality and other basic values is out of the scope of our
362 paper, we can deal with another argument against *RFV* that is rather closely
363 related to the problem of moral responsibility.

364 Namely, *RFV* seems to imply that epistemic obligations are identical with
365 or, at least, not independent of moral obligations which means, in turn, that
366 anyone who accepts *RFV* and would like to believe in epistemic norms is
367 forced to believe in moral responsibility no matter which crazy theory about
368 conditions of moral responsibility turns out to be true. For example, if Derk
369 Pereboom’s analysis on the conditions of moral responsibility is correct, moral
370 responsibility needs not only agent-causation (which, according to Pereboom,
371 may be a logically incoherent concept), but either systematic breaking of the
372 laws of nature or inexplicable harmony between micro-physical statistical
373 laws and the free decisions of the agents (see Pereboom 2001). For sure,
374 believing that these conditions are met in reality would be a high price to
375 pay for holding *RFV*. Insofar as the price is so high, it seems to be not only
376 irrelevant, but weird that the proponent of *RFV* can and even should rationally
377 defend believing in moral responsibility and its monstrous metaphysics by
378 moving the battlefield from metaphysics to metaethics. After all, forming
379 rational beliefs and fulfilling epistemic norms aim at the truth, and it is not
380 too probable that this way of belief-formation leads us to true beliefs.⁵

381 To be clear, *RFV* does not imply that epistemic obligations as such depend
382 on (or are identical with) moral obligations. *RFV* does not exclude that they
383 are totally unrelated to each other. What *RFV* implies is only that an epistemic
384 obligation has to be appropriately related to some moral obligations in order
385 to be true that agents have to avoid violating it with all their strength and
386 if someone fails to observe an epistemic obligation which does not relate
387 appropriately to any moral obligation, she should not feel remorse. In other
388 words, in themselves, epistemic obligations do not have sufficient normative
389 weight to constitute all things considered obligations. So, if one both accepts
390 *RFV* and rejects moral responsibility based on her evidence-basis, she cannot
391 rationally believe that there are moral and all things considered obligations,

5 We would like to thank an anonymous reviewer for pointing out this possible objection.

392 but she still can rationally think that there are (weak) epistemic obligations.
 393 What she cannot rationally believe is that neglecting any epistemic obligation
 394 cannot be OK in the same way as neglecting the prudential obligation not to
 395 ridicule oneself. That is, even if one accepts **RFV** and, for instance, Pereboom's
 396 assessment of the evidence about free will and moral responsibility, she can
 397 rationally deny the existence of (a metaphysically rather extreme kind of)
 398 free will, moral responsibility, moral obligations, and all things considered
 399 obligations. The only thing that she cannot rationally maintain without re-
 400 jecting **RFV** is the idea that anyone should suspend belief in those things to
 401 avoid committing an unacceptable normative error that cannot be shrugged
 402 off. In other words, if there is a proponent of **RFV** who tries to heroically
 403 find the truth no matter the cost and finds that her evidence-basis strongly
 404 indicates the non-existence of moral responsibility, she can rationally believe
 405 that she has an epistemic (or even prudential) obligation to deny the existence
 406 of moral responsibility, but she cannot rationally think either that she has an
 407 all things considered obligation to reject moral responsibility or that anyone
 408 has an all things considered obligation to try to find the truth no matter the
 409 cost.

410 Nonetheless, none of this undermines our point that the proponent of
 411 **RFV**, if she wants to defend the belief in moral responsibility, can move the
 412 battlefield from metaphysics to metaethics, and the latter seems to be much
 413 more advantageous for her, especially if she also holds that the sufficient
 414 conditions of moral responsibility are metaphysically rather demanding. The
 415 more demanding these conditions are, the less plausible is the claim that the
 416 existence of moral responsibility is obvious and/or probable in the light of
 417 the given evidence, so moving the battlefield from metaphysics to metaethics
 418 provides more strategic advantage.

419 Note, parenthetically, that one can construct a modified version of the
 420 **PRIMACY ARGUMENT** even if both **MPT** and **RIT** are untenable. One need
 421 not appeal to morality (or moral responsibility) at all. Anyone who rationally
 422 believes in all things considered obligations has the epistemic right to sustain
 423 belief in a specific kind of responsibility. As the first premise of the **PRIMACY**
 424 **ARGUMENT** says, no agent can rationally think she has an all things consid-
 425 ered obligation to suspend the belief that she fulfills the necessary conditions
 426 of having all things considered obligations. In other words, someone who ra-
 427 tionally attributes all things considered obligations to herself must also accept
 428 that she fulfills the necessary conditions of having all things considered obliga-
 429 tions. Since having strong obligations is one of those necessary conditions, the

430 agent in question must also accept that she fulfills the necessary conditions
431 of having strong obligations. It is reasonable to say that being responsible “in
432 a strong sense” requires fulfilling all necessary conditions for bearing strong
433 obligations, so anyone who rationally attributes all things considered obliga-
434 tions to herself can rationally attribute “strong responsibility” to herself as
435 well. It seems that this argument for “strong responsibility” can be threatened
436 only by arguments against the existence of all things considered obligations.

437 To sum up the **PRIMACY ARGUMENT**, anyone who rationally accepts **RFV**
438 can rationally maintain her belief in moral responsibility even if she does not
439 have sufficient direct evidence that anyone fulfills the metaphysical conditions
440 of being morally responsible. Until the agnostic refutes **MPT** or **RIT**, one can
441 rationally resist the agnostic challenge.

446 **6 Free Will, Fairness, and Others**

443 Various questions could be raised about our argument. We will look at three.
444 First, one might ask how the dialectic is related to free will. We claimed that
445 someone who rationally believes in **RFV** does not have to suspend her belief
446 in moral responsibility even if she lacks direct evidence for it. Could **RFV** be
447 used to defend belief in free will as well?

448 The answer to this question depends on one’s conception of free will. There
449 are two basic approaches in the literature. According to the first, having free
450 will means fulfilling a subset of conditions that guarantee necessary control
451 over one’s morally relevant actions (Clarke 1992). The present argument
452 obviously extends to the defense of free will conceived this way. If someone
453 rationally accepts **RFV** and also rationally thinks that she fulfills all necessary
454 conditions for being morally responsible, then she cannot rationally believe
455 that she fails to fulfill a subset of those conditions, namely, those that are
456 necessary for control. So our argument supports belief in free will for those
457 who rationally believe **RFV** and identify having free will with fulfilling a
458 subset of necessary control conditions for being morally responsible.

459 However, there is another prevalent conception, according to which free
460 will is the ability to do otherwise (van Inwagen 1983). Our argument can be
461 extended to this case as well, but only if one rationally upholds that the ability
462 to do otherwise is a necessary precondition of being morally responsible. Given
463 strong evidence that moral responsibility depends on free will of the second
464 sort, then rational belief in **RFV** (together with the evidence in question) can
465 ground rational belief in the existence of free will. And if the proponent of

466 RFV has sufficient evidence that moral responsibility has further metaphysical
467 conditions, she can also rationally believe that she fulfills all those further
468 conditions, regardless of how demanding they are metaphysically.

469 These possible extensions of the **PRIMACY ARGUMENT** are especially sig-
470 nificant if one takes into account that many philosophers and scientists insist
471 that there is no sufficient scientific or other evidence for macro-level psycho-
472 logical indeterminism (which is a precondition of libertarian free will) or
473 the presence of compatibilist-friendly causal determinism in the brain. In
474 light of the possible extensions of the **PRIMACY ARGUMENT**, proponents of
475 RFV can rationally believe in responsibility-relevant free will (of either the
476 incompatibilist or compatibilist sort) even in the absence of sufficient direct
477 scientific or other evidence.

478 This last point regarding the absence of evidence leads us on to a potential
479 objection implied by Scott Sehon. He emphasizes that we treat responsible
480 and irresponsible agents very differently. If, for example, someone pushes
481 another person into the traffic, we treat her act very differently depending on
482 whether she was or was not responsible. If she was, then her act “certainly
483 looks incredibly reprehensible and maybe even the stuff of an attempted
484 murder charge” (Sehon 2013, 369). But if we know that the pusher is not
485 responsible, we would not call her action “reprehensible” and would not
486 make her face serious charges. Sehon adds that “[it] would be manifestly
487 unfair to regard the agent as responsible if our degree of certainty on the
488 matter is quite low” (2013, 36). One could extend this point and argue that if
489 we lack strong direct evidence for moral responsibility, then, out of fairness,
490 we should suspend belief about whether anyone is ever responsible in a way
491 that would render retribution justified.

492 Proponents of the **PRIMACY ARGUMENT** evidently disagree, as their sup-
493 posedly rational belief in RFV makes them rational in holding that they can
494 be morally responsible for their actions. It is important that Sehon brings up
495 this issue in terms of fairness. The obligation to be fair with others is naturally
496 understood to be a moral obligation, and hence a strong one—according to
497 MPT. Those who uphold RFV will see the situation as follows. The obligation
498 to be fair can only be attributed to morally responsible persons. If nobody is
499 morally responsible, then the strong obligation to be fair cannot be attributed
500 to anybody. And if that is the case, then nobody has to care about being fair
501 to anybody. So if the proponent of RFV takes Sehon’s exhortation to be fair
502 seriously, and if she thinks she has to care about it, then, in the light of MPT,
503 she incidentally attributes a strong obligation to herself. As a result, she im-

504 plies that she fulfills all the conditions of having strong obligations, including
505 having moral responsibility. That is, for proponents of RFV, Sehon's point can
506 only have force if they take themselves to be morally responsible. They would
507 need to assume, first, that they are morally responsible, and, second, they
508 would have to suspend judgment about the existence of moral responsibility
509 because of that very assumption—which seems incoherent. Thus, the argu-
510 ment that insufficient direct evidence for moral responsibility should make us
511 suspend our belief in moral responsibility because it might lead to the unfair
512 treatment of others makes no sense to those who hold RFV to be true.

513 The question of being fair to others and taking them to be morally respon-
514 sible brings us to the crucial issue of the circle of agents whom one might
515 attribute moral responsibility to, on the basis of the PRIMACY ARGUMENT.
516 This is a crucial issue, as it might very well be the case that the individuals
517 who accept RFV can attribute moral responsibility only to themselves and
518 not to anyone else. This is because (PA1) takes only the agent's own perspec-
519 tive into account. The agent is considering her own obligations and moral
520 responsibility, and the reason why she doesn't have to become an agnostic is
521 that, were she to take agnosticism as a strong obligation, she would thereby
522 attribute moral responsibility to *herself*. (As we have indicated, she might even
523 go on to attribute free will as well.) But the incoherence would arise only in
524 her own case, so the PRIMACY ARGUMENT's conclusion applies only to her:
525 she is free to go on believing that she, for one, is morally responsible. And
526 clearly, she can believe in the existence of moral responsibility on the basis
527 on *that*, since moral responsibility exists even if only one agent has it.

528 Extending this rather small circle of responsible agents might look unrea-
529 sonable or unfair indeed. However, there could be ways to do it. Remember
530 that the agent in question reasonably believes in her own responsibility. If she
531 considers agents who seem to be like her in every relevant respect, she may
532 take them to be morally responsible as well. Nevertheless, the reasonableness
533 of this move depends on two crucial factors. First, the agent must have a
534 rationally held theory of what the relevant respects are. Second, were she to
535 deem morally responsible any agent other than herself, her judgment that that
536 person is similar to herself in every relevant respect must also be rational. This
537 means that reasonably extending the circle of morally responsible agents to
538 others is logically possible, but could be difficult in practice. Fortunately, there
539 might be an easier way. It seems reasonable to think that all fully developed
540 human beings have the same metaphysical structure. Insofar as this assump-
541 tion is reasonable, a proponent of the Primacy Argument can extend the circle

of morally responsible agents to all fully developed human adults who fulfill the non-metaphysical and empirically verifiable conditions of moral responsibility, whether or not she can identify the precise metaphysical conditions for being morally responsible.

Note that extending the circle of responsibility poses a challenge not only with regard to other agents, but also with regard to the agent who can rationally believe in moral responsibility based on *RFV* and *PA*. This is because the *PRIMACY ARGUMENT* does not imply that the agent is morally responsible all the time. It only permits the agent to believe that she is morally responsible in her present state. Nevertheless, what we have said previously about the possibilities of extending the circle of morally responsible agents can also be used to extend this temporal limitation. This means that if an earlier or later state of the agent seems to be similar in every relevant respect to her present state, then she may take it that she was or is going to be morally responsible at those times. However, it might not be clear in every case that these conditions are fulfilled. Therefore, our argument is compatible with accepting that even though we are reasonable in thinking that we are morally responsible some of the time, we could be also reasonable in thinking that we are not responsible at other times, or thinking that we should be agnostics about the question whether we are morally responsible in certain situations.*

László Bernáth

Research Centre for the Humanities, Budapest

bernath.laszlo@abtk.hu

* We would like to express our indebtedness to the anonymous referees of this paper, their insightful comments helped us developing our position a great deal. We also thank Dániel Kodaj for his suggestions and the stimulating exchanges we had on many occasions. László Bernáth was supported by the János Bolyai Research Scholarship of the Hungarian Academy of Sciences (grant no. BO/00432/18/2), the OTKA (Hungarian Scientific Research Fund by the National Research Development and Innovation Office) Postdoctoral Excellence Programme (grant no. PD131998), and another OTKA research grant (grant no. K132911). Tamás Paár was supported by the ÚNKP-19-3-III-PPKE-24 New National Excellence Program of the Ministry for Innovation and Technology. His research was supported by the University of Oxford project “New Horizons for Science and Religion in Central and Eastern Europe” funded by the John Templeton Foundation. (The opinion expressed in the publication are those of the author(s) and do not necessarily reflect the view of the John Templeton Foundation.) We owe special thanks to the organizers and the committee of the Oxford University Essay Competition for Central and Eastern European Perspectives on Science, Theology and Philosophy in which this paper has been awarded first prize.

Tamás Paár

 0000-0003-0585-6388

Central European University, Budapest

paar.tamas@gmail.com

References

- BYRD, Jeremy. 2010. "Agnosticism about Moral Responsibility." *Canadian Journal of Philosophy* 40(3): 411–432, doi:[10.1080/00455091.2010.10716729](https://doi.org/10.1080/00455091.2010.10716729).
- . 2021. "What Should We Believe About Free Will?" *Erkenntnis* 86(3): 505–522, doi:[10.1007/s10670-019-00116-3](https://doi.org/10.1007/s10670-019-00116-3).
- CLARKE, Randolph. 1992. "Free Will and the Conditions of Moral Responsibility." *Philosophical Studies* 66(1): 53–72, doi:[10.1007/bf00668395](https://doi.org/10.1007/bf00668395).
- COFFMAN, E. J. 2016. "Incompatibilist Commitment and Moral Self-Knowledge: The Epistemology of Libertarianism." in *Philosophical Issues 26: Knowledge and Mind*, edited by Christoph KELP and Jack C. LYONS, pp. 78–98. Malden, Massachusetts: Wiley-Blackwell, doi:[10.1111/phis.12066](https://doi.org/10.1111/phis.12066).
- FELDMAN, Richard H. 1988. "Epistemic Obligations." in *Philosophical Perspectives 2: Epistemology*, edited by James E. TOMBERLIN, pp. 235–256. Atascadero, California: Ridgeview Publishing Co., doi:[10.2307/2214076](https://doi.org/10.2307/2214076).
- GUILLOIN, Jean-Baptiste. 2014. "Van Inwagen on Introspected Freedom." *Philosophical Studies* 168(3): 645–663, doi:[10.1007/s11098-013-0159-x](https://doi.org/10.1007/s11098-013-0159-x).
- VAN INWAGEN, Peter. 1983. *An Essay on Free Will*. Oxford: Oxford University Press.
- KEARNS, Stephen. 2015. "Free Will Agnosticism." *Noûs* 49(2): 235–252, doi:[10.1111/nous.12032](https://doi.org/10.1111/nous.12032).
- LATHAM, Andrew James. 2019. "The Conceptual Impossibility of Free Will Error Theory." *European Journal of Analytic Philosophy* 15(2): 99–120, doi:[10.31820/eja.p.15.2.5](https://doi.org/10.31820/eja.p.15.2.5).
- LEVY, Neil. 2011. *Hard Luck: How Luck Undermines Free Will and Moral Responsibility*. Oxford: Oxford University Press, doi:[10.1093/acprof:oso/9780199601387.001.0001](https://doi.org/10.1093/acprof:oso/9780199601387.001.0001).
- LOCKIE, Robert. 2018. *Free Will and Epistemology: A Defence of the Transcendental Argument for Freedom*. London: Bloomsbury Academic, doi:[10.5040/9781350044180](https://doi.org/10.5040/9781350044180).
- MELE, Alfred R. 2019. "Causalism: On Action Explanation and Causal Deviance." in *Explanation in Action Theory and Historiography. Causal and Teleological Approaches*, edited by Gunnar SCHUMANN, pp. 45–58. London: Routledge, doi:[10.4324/9780429506048-2](https://doi.org/10.4324/9780429506048-2).
- PEREBOOM, Derk. 2001. *Living Without Free Will*. Cambridge: Cambridge University Press, doi:[10.1017/cbo9780511498824](https://doi.org/10.1017/cbo9780511498824).
- ROSS, William David. 1930. *The Right and the Good*. Oxford: Oxford University Press.

- 603 RUSSELL, Bruce. 2001. "Epistemic and Moral Duty." in *Knowledge, Truth, and Duty:*
604 *Essays on Epistemic Justification, Responsibility, and Virtue*, edited by Matthias
605 STEUP, pp. 34–48. Oxford: Oxford University Press, doi:10.1093/0195128923.003.0
606 003.
- 607 SEARLE, John R. 1978. "Prima Facie Obligations." in *Practical Reasoning*, edited by
608 Joseph RAZ, pp. 81–90. Oxford: Oxford University Press. Reprinted in van Straaten
609 (1980, 238–259).
- 610 SEHON, Scott R. 2013. "Epistemic Issues in the Free Will Debates: Can we Know when
611 We Are Free?" *Philosophical Studies* 166(2): 363–380, doi:10.1007/s11098-012-
612 0044-z.
- 613 —. 2016. *Free Will and Action Explanation: A Non-Causal, Compatibilist Account*. Oxford:
614 Oxford University Press, doi:10.1093/acprof:oso/9780198758495.001.0001.
- 615 SINNOTT-ARMSTRONG, Walter. 2008. "Framing Moral Intuitions." in *Moral Psychology,*
616 *Volume 2. The Cognitive Science of Morality: Intuition and Diversity*, edited by
617 Walter SINNOTT-ARMSTRONG, pp. 47–76. Cambridge, Massachusetts: The MIT
618 Press, doi:10.7551/mitpress/7573.003.0004.
- 619 SMILANSKY, Saul. 1994. "The Ethical Advantages of Hard Determinism." *Philosophy*
620 *and Phenomenological Research* 54(2): 355–363, doi:10.2307/2108494.
- 621 VAN STRAATEN, Zak, ed. 1980. *Philosophical Subjects: Essays Presented to P.F. Strawson*.
622 Oxford: Oxford University Press.
- 623 STRAWSON, Galen. 1994. "The Impossibility of Moral Responsibility." *Philosophical*
624 *Studies* 75(1/2): 5–24. Reprinted as "The Impossibility of Ultimate Moral Respon-
625 sibility" in Strawson, G. (2008, 319–336), doi:10.1007/bf00989879.
- 626 —. 2008. *Real Materialism, and Other Essays*. Oxford: Oxford University Press, doi:10
627 .1093/acprof:oso/9780199267422.001.0001.
- 628 STRAWSON, Peter Frederick. 1969. *Meaning and Truth*. Oxford: Oxford University Press.
629 An Inaugural Lecture Delivered Before the University of Oxford on 5 November
630 1969. Reprinted in Strawson, P.F. (1971, 170–189) and in Strawson, P.F. (2004,
631 131–146).
- 632 —. 1971. *Logico-Linguistic Papers*. London: Methuen & Co. Reprinted as Strawson, P.F.
633 (2004), doi:10.4324/9781351153607.
- 634 —. 2004. *Logico-Linguistic Papers*. 2nd ed. London: Routledge, doi:10.4324/97813152
635 50250.
- 636 TRAKAKIS, Nick N. 2007. "Whither Morality in a Hard Determinist World?" *Sorites*
637 19(1): 14–40, http://lorenzopena.es/Sorites/Issue_19/sorites19.pdf.
- 638 WALLER, Bruce N. 2004. "Virtue Unrewarded: Morality Without Moral Responsibility."
639 *Philosophia: Philosophical Quarterly of Israel* 31(3/4): 427–447, doi:10.1007/bf0238
640 5194.
- 641 WILLIAMS, Bernard Arthur Owen. 1985. *Ethics and the Limits of Philosophy*. London:
642 Fontana Press.

- 643 —. 2006. *Ethics and the Limits of Philosophy*. London: Routledge. With a commentary
644 on the text by A.W. Moore; original publication: Williams (1985), doi:[10.4324/97](https://doi.org/10.4324/9780203969847)
645 [80203969847](https://doi.org/10.4324/9780203969847).
- 646 WITTGENSTEIN, Ludwig. 1965. "A Lecture on Ethics." *The Philosophical Review* 74(1):
647 3–12, doi:[10.2307/2183526](https://doi.org/10.2307/2183526).

PROOF

PROOF

Self-Knowledge and Interpersonal Reasoning

BENJAMIN WINOKUR

Many philosophers contend that we often possess “privileged” and “peculiar” self-knowledge of our mental states. Self-knowledge is privileged insofar as it is systematically more secure than the knowledge that others have of one’s propositional attitudes, and it is peculiar insofar as it is systematically obtained in a way that is only suited for delivering self-knowledge. Focusing on privileged and peculiar self-knowledge of propositional attitudes like beliefs, I offer an account of its instrumental value. On my account, privileged and peculiar self-knowledge of one’s propositional attitudes enables one to be a more efficient and reliable interpersonal reasoner.

Self-knowledge of one’s current mental states often seems interesting—if not outright puzzling—for at least two reasons. First, such self-knowledge often seems to be *privileged*, for it seems to be systematically (though not universally) more secure than the knowledge one has of others’ mental states. Second, it often seems to be *peculiar*, for it seems to be systematically (though, again, not universally) obtained in a way that is only suited for delivering self-knowledge, hence, *not* by whatever means enable one to acquire knowledge of *other* minds (Byrne 2018, 4–9). The standard project in contemporary theorizing about self-knowledge is to vindicate these appearances by unearthing the special security and sources of self-knowledge. However, others have argued that we do not actually possess any privileged and peculiar self-knowledge (hereafter “PPSK”), at least when it comes to self-knowledge of propositional attitudes like belief (Gopnik 1993; Carruthers 2011; Cassam 2014). These PPSK-skeptics typically understand self-knowledge and other-knowledge of propositional attitudes as on a par in terms of their security, source, or both.

In reply, some PPSK-realists have offered competing interpretations of the putative evidence against realism about PPSK of propositional attitudes (Parent 2017; Keeling 2019a, 2019b; Andreotta 2021; Valaris 2018), while

678 others have pushed back against the non-privileged and non-peculiar accounts
679 of self-knowledge that are favoured by many skeptics (Coliva 2016; Keeling
680 2018; Marcus and Schwenkler 2019; Andreotta 2022). The stakes of these
681 debates are hard to grasp if we are unsure “what, if anything, of value we fail
682 to possess if these skeptics are right” (Peterson 2021, 365). For this reason, I
683 will argue that PPSK of one’s propositional attitudes—chiefly, our beliefs—
684 is instrumentally valuable for the efficiency and reliability of a widespread
685 activity in our social-epistemic lives, that of *interpersonal reasoning*. Some
686 readers may interpret my arguments as providing further support for PPSK-
687 realism if they believe that interpersonal reasoning is in fact a highly efficient
688 and reliable activity in our actual lives. Other readers might reach the more
689 modest conclusion that interpersonal reasoning is a more effective enterprise
690 *to the extent that we possess PPSK*, whether or not we really possess PPSK, and
691 hence whether or not interpersonal reasoning is a particularly efficient and
692 reliable activity for us to undertake. Either way, the significance of debates
693 between PPSK-skeptics and PPSK-realists can be better appreciated in light
694 of what follows.

695 Here is the layout for my paper. In section 1 I draw initial inspiration from
696 two earlier accounts of PPSK’s instrumental value. The first, due to Sydney
697 Shoemaker (1988, 1996), concludes that social cooperation in general requires
698 each of us to possess PPSK of many of our propositional attitudes. The second,
699 due to Charles Siewert (2003), concludes that PPSK is indispensable to social
700 cooperation *whenever this depends on justifying one’s actions to others*. Justify-
701 ing one’s actions to others can be one way of reasoning with others, that is,
702 reasoning interpersonally. But it is *only* one way of reasoning interpersonally.
703 I thus consider, in section 2, whether *all* interpersonal reasoning might benefit
704 from PPSK. My argument is that PPSK does indeed play beneficial roles in all
705 interpersonal reasoning. In section 3 I address objections to my account. In
706 section 4 I consider another recent account of PPSK’s instrumental value, one
707 that emphasizes its role in our capacity for “epistemic control,” and I show
708 how my account complements that account. In section 5 I conclude.

7091 Cooperation and Privileged, Peculiar Self-Knowledge

710 Is PPSK instrumentally valuable? Some philosophers have argued that it is.¹
 711 Indeed, some have argued that it is instrumentally *indispensable*. Here is
 712 Shoemaker, who writes of self-knowledge by “self-acquaintance” instead of
 713 privileged and peculiar self-knowledge:²

714 When one is engaged in a cooperative endeavor with another, it
 715 is essential to the efficient pursuit of the shared goal that one be
 716 able to communicate to the other information about one’s beliefs,
 717 desires and intentions [...] When in such circumstances one conveys
 718 one’s beliefs to another, this is not merely for the purpose
 719 of conveying what one takes to be information about the world,
 720 namely, the contents of the beliefs; it is also for the purpose of
 721 giving him information about oneself which will assist him in
 722 predicting one’s behavior and so in coordinating his own behavior
 723 with it, and also to enable him to correct those of one’s beliefs he
 724 knows to be mistaken [...] And here the utility of self-knowledge
 725 depends crucially on its being acquired by self-acquaintance; if
 726 I had to figure out from my behavior what my beliefs, goals, in-
 727 tentions, etc. are, then in most cases it would be more efficient
 728 for others to figure this out for themselves than to wait for me to
 729 figure it out and then tell them about it. (1988, 185–186)

730 Shoemaker argues that PPSK is indispensable for efficiently cooperating with
 731 other human beings. For, if others could know one’s mind in the same way
 732 and as reliably as one knows one’s own mind, one would be far less efficient
 733 at soliciting others’ cooperation. This is because it would just as often be
 734 up to others to figure out one’s mind, and to decide on this basis whether
 735 cooperation was worthwhile. As a result, one would frequently fail to solicit
 736 others’ cooperation of one’s own accord.

1 Peterson (2021, 365) thinks that the question of PPSK’s value has been ignored by epistemologists working on self-knowledge. While I myself hope to contribute an answer to this question, I think that this assessment of the extant literature is somewhat exaggerated given the views that I discuss in this section, among others (see, e.g., Burge 1996; Nguyen 2015; Sorgiovanni 2019; Winokur 2021a, 2021b). Peterson’s own account of PPSK’s instrumental value is discussed in section 4. He also discusses the potential *intrinsic* value of PPSK, a topic that I do not broach here.

2 These, I submit, are just notational variants.

737 Reflecting on Shoemaker’s argument, Siewert wonders whether rational
 738 animals “could engage in cooperation and assistance-seeking behaviour, even
 739 by generally acting in an attitude-revealing fashion, without representing
 740 their own minds to themselves” (2003, 139). In a different idiom: couldn’t
 741 there be creatures that are exceptionally adept at *expressing*—i.e., showing,
 742 manifesting, displaying—their attitudes to their fellow creatures without also
 743 possessing *PPSK of the attitudes expressed*, and couldn’t this enable equally
 744 efficient cooperation?³ Contra Shoemaker, Siewert supposes that there could
 745 be such creatures. Still, he is optimistic about a nearby argument:

746 For whether or not there can be social animals that act in a usefully
 747 self-revealing fashion while oblivious to their own psychologies,
 748 they could not engage in the practice of *justifying* such acts, with-
 749 out being able to represent, in their justifications, relevant facts
 750 about their own desires and beliefs [...] Now, if the reasons we
 751 would offer did not have us acting in ways revealing our actual
 752 beliefs and desires to others, we would be much less effective in
 753 securing others’ cooperation and assistance in the satisfaction of
 754 our desires than we in fact are. (Siewert 2003, 139)

755 On this argument, it is not that efficient cooperation always requires *PPSK*.
 756 Rather, such self-knowledge is required for cooperation *whenever such cooper-*
 757 *ation also depends on justifying one’s actions to one’s would-be cooperators*. For,
 758 lacking *PPSK*, our actions would often fail to cohere with the attitudes that
 759 we self-ascribe. In turn, we would be worse at justifying our actions because
 760 we would be worse at appealing to the actual beliefs, desires, and intentions
 761 that underwrite them. These inconsistencies might be noticed by others, and
 762 this might diminish their trust in us.

763 More recently, Jon Greco has written that:

764 *Of course* thinking about one’s first-order mental states is essen-
 765 tial to activities involving coordination and cooperation [...] In
 766 particular, *giving one’s reasons*, both epistemic and practical, is es-
 767 sential to various activities in which one must defend one’s beliefs
 768 and actions, and having a grasp on such mental states oneself is
 769 essential to reporting them to others. (2019, 52)⁴

3 This is my gloss on Siewert’s argument. Like Bar-On (2004), I use “express” here to denote actions that express mental states, though I denote another sense of expression in section 3.

4 See Müller (2019, 6) for a similar view.

770 Like Siewert, Greco claims that self-knowledge is essential for justifying one's
771 actions (and beliefs) to others. However, Siewert argues that PPSK is indispens-
772 able to our widespread success in these matters, whereas Greco claims that
773 “this kind of metacognitive activity can tolerate the same fallibility that we
774 experience in cognition generally” (2019, 53).⁵ He thus denies the importance
775 of *privileged* self-knowledge (and is silent about *peculiar* self-knowledge). But
776 he does not consider Siewert's argument, and so it is hard to know whether
777 his position would change upon further reflection.

778 This difference between Greco and Siewert set aside, notice that they both
779 focus on a certain kind of interpersonal reasoning. Here, ‘interpersonal rea-
780 soning’ denotes exchanges of assertions between interlocutors, or exchanges
781 of questions and assertions, toward a discursive end. For instance, one might
782 reason interpersonally in order to acquire new rational attitudes, or to subject
783 one's already-held attitudes to the scrutiny of other agents, or to persuade
784 other agents to adopt one's already-held attitudes. I say that Siewert and Greco
785 are focused on a certain kind of interpersonal reasoning because they only
786 focus on cases in which agents reason interpersonally about one another's
787 actions or attitudes. In other words, neither philosopher focuses on cases
788 in which agents aim to justify “agent-neutral” propositions to one another,
789 these being propositions whose contents do not refer to any particular agent's
790 actions or attitudes. One such proposition is:

791 Runaway climate change is a worsening phenomenon.

792 It is to be contrasted with the sorts of propositions that Siewert and Greco
793 focus on, namely, “agent-specific” propositions like:

794 I should continue to be vigilant about my fossil fuel consumption.

5 Greco also ventures a response to the possibility of efficiently cooperative animals lacking self-knowledge: “One might object that non-human animals are also social in a sense that implies coordination and cooperation, and they manage their social lives without citing their mental states in explanations to themselves or their cohorts. But this objection misses the point that human social agency is also rational agency. It involves rationalizing one's thoughts and actions by means of giving one's reasons—i.e., overtly giving one's reasons—to oneself and to others” (Greco 2019, 53). This too is reminiscent of Siewert's view. But while Greco denies that cooperation among non-human animals involves rational agency, Siewert thinks that non-human animals could count as rational agents in so cooperating.

795 The latter proposition, but not the former, requires agents to provide self-
 796 referential information, this being information that justifies *the agent herself*
 797 to act in such-and-such a way or have such-and-such an attitude. Such infor-
 798 mation will naturally include “relevant facts about their own desires and be-
 799 liefs” (Siewert 2003, 139) whereas reasoning about agent-neutral propositions
 800 simply requires providing first-order evidence about the agent-independent
 801 world, e.g., evidence of rapidly melting arctic ice. But because reasoning about
 802 either sort of proposition can be conducted interpersonally, Siewert and Greco
 803 will have shown—at most—that self-knowledge matters for interpersonal
 804 reasoning about agent-specific matters (whether such self-knowledge is privi-
 805 leged and peculiar as Siewert claims, or not, as Greco claims). I note this here
 806 because I will argue in section 2 that PPSK plays a role in both agent-specific
 807 and agent-neutral interpersonal reasoning.

808 Before I get there, I want to make two preliminary points. First, the reader
 809 may have wondered whether that Shoemaker’s and Siewert’s arguments estab-
 810 lish what they purport to since, on close inspection, they seem to emphasize
 811 the importance of privileged access but not, in addition, peculiar access. This
 812 is because each argument insists that the special security of agents’ self-
 813 knowledge is what facilitates cooperation with other people, and yet this does
 814 not obviously entail that agents must exploit a peculiar means of achieving
 815 such security. In what follows I will provide arguments for the importance of
 816 peculiar self-knowledge as well, thereby going beyond the arguments consid-
 817 ered thus far.

818 Second, it should be noted that some philosophers deny that interpersonal
 819 reasoning of *any* kind (i.e., whether about agent-neutral or agent-specific
 820 propositions) requires self-knowledge of *any* kind (i.e., whether privileged
 821 and peculiar, non-privileged and non-peculiar, or any other combination).
 822 For example, Robert Brandom writes that there is “nothing incoherent in
 823 descriptions of communities of judging and perceiving agents, attributing
 824 and undertaking propositionally contentful commitments, giving and asking
 825 for reasons, who do not yet have available the expressive resources *I* pro-
 826 vides” (1994, 559). If these communities lack articulate use of the first-person
 827 singular, then they cannot self-ascribe and hence self-know their attitudes.⁶

6 See also Stribos and de Bruin (2012). The importance of this claim depends on assuming that self-knowledge requires linguistically articulate self-ascriptive thought, and some friends of “tacit” self-knowledge might dispute this (e.g., Boyle 2011, 2019). Alternatively, it could be granted that there is such a thing as self-consciousness that does not involve linguistically articulate self-ascriptive thought (cf. Musholt 2015, chap. 4). Even if this is a tenable view, I am focusing

828 Similarly, Ladislav Koreň claims that we can reason interpersonally by exer-
829 cising a “practical competence” with linguistic devices like “no,” “but,” and
830 “so,” thus manifesting a “sensitivity” to rational connections between claims
831 without having “metarepresentational” beliefs about the rational connec-
832 tions between one’s own attitudes or one’s interlocutor’s attitudes (2023, 853).
833 Finally, Annalisa Coliva offers the following thought experiment:

834 Take a subject who is able to judge that P, give evidence in favour
835 of it and withdraw from it if required and, therefore, has the first-
836 order belief that P based on judgement. Suppose you ask her “Do
837 you believe that P?” and she is unable to answer. You conclude
838 that she does not have the concept of belief. (2016, 191)

839 This is a situation in which one interlocutor reasons interpersonally while, *ex*
840 *hypothesi*, lacking the conceptual wherewithal to self-ascribe the attitudes that
841 her assertions express. Coliva adds that any such agent will at least possess
842 “the *ability* to differentiate between, for instance, believing P and P’s being the
843 case, by being *sensitive to the fact that* her point of view may be challenged [...]”
844 (2016, 192, emphasis mine). On my reading, the emphasized terms suggest
845 that such an agent utilizes *pre-metarepresentational* capacities in the service of
846 interpersonal reasoning; these abilities and sensitivities enable her to reason
847 with others without forming second-order beliefs about her first-order beliefs
848 or her interlocutor’s first-order beliefs.

849 These philosophers clearly reject Greco’s claim that “thinking about one’s
850 first-order mental states is essential to activities involving coordination and
851 cooperation,” given that interpersonal reasoning is itself a coordinated and co-
852 operative endeavour. But do they extend this rejection as far as to deny that in-
853 terpersonal reasoning *with an aim to justifying one’s own actions and attitudes*
854 requires self-knowledge or, at the very least, some form of self-representation
855 like a self-belief? As Steven Levine makes clear in a response to Brandom, it
856 is hard to see how they could cogently deny this. Levine begins by acknowl-
857 edging the possibility of agents who reason interpersonally insofar as the
858 assertions at issue are first-order assertions of the form “that-P,” these being
859 expressions of agent-neutral propositions in the sense described above. As
860 regards assertions of these propositions, “the performer can justify the state-
861 ment without explicitly claiming that it is he who is justifying the statement

on what epistemologists in this area ordinarily focus on, i.e., *explicit* self-knowledge involving linguistically articulate self-ascriptive thought (*pace* also those who view self-knowledge as an *ability*—cf. Campbell 2018).

862 [...] because this assertion concerns an objective state of affairs that can be
863 justified by agent-neutral reasons” (Levine 2009, 111). However:

864 [...] is this the case when the assertion that is being challenged
865 concerns an agent’s *own* action or perception? Here what is being
866 challenged is, for example, one’s entitlement to perform an action
867 or one’s entitlement to claim that one’s perception is veridical.
868 In either case, the justificatory reasons offered cannot be agent-
869 neutral in the way that reasons justifying the assertion “that-P”
870 are. (Levine 2009, 111)

871 So Levine is in league with Siewert and Greco in arguing that, when one’s own
872 actions are challenged by an interlocutor, one cannot merely avail oneself of
873 agent-neutral reasons. Instead, one must avail oneself of agent-specific rea-
874 sons, which will include facts about one’s own psychology. The only question
875 is whether Levine would side with Siewert in understanding these exchanges
876 as requiring PPSK on the part of whoever seeks to justify her own attitudes,
877 or with Greco in denying any indispensable role for such epistemically high-
878 grade self-knowledge.

879 As aforementioned, I will soon argue that PPSK plays important roles in
880 interpersonal reasoning about both agent-specific and agent-neutral proposi-
881 tions. But how can I be headed in this direction, having just traced a dialectic
882 that only acknowledges a role for self-knowledge in interpersonally defending
883 agent-specific propositions about one’s own actions or perceptions? In other
884 words, if it is conceded to Brandom and others that agents can reason inter-
885 personally about agent-*neutral* propositions without so much as a capacity
886 for self-belief, then isn’t it foolish to contend that PPSK—let alone any other
887 sort of self-knowledge—matters for such activity? Fortunately, there is no
888 real problem here. My argument will be that PPSK contributes to interper-
889 sonal reasoning *for agents who in fact possess the capacity for representing*
890 *themselves and their beliefs in higher-order thought*. This focus allows me to
891 grant Brandom, Koreň, and Coliva their contention that some agents can
892 reason interpersonally despite lacking this metarepresentational capacity.⁷

7 There are other ways to dispute the indispensability of self-knowledge for interpersonal reasoning. For example, Roelofs (2017) argues that no such knowledge is required by interpersonal reasoners who are “evidentially unified” with and “cognitively vulnerable” to one another. Evidentially unified agents are automatically attuned to one another’s evidence without having to explicitly share it, while cognitively vulnerable agents can rationally cause changes in one another’s minds through cognizing this unified evidence (they can induce such changes as *basic actions*).

893 What I will argue is that agents who *do* possess this capacity, such as most
 894 cognitively developed adult human beings, are systematically vulnerable to
 895 certain deficiencies in interpersonal reasoning if they lack **PPSK**.

892 **Interpersonal Reasoning and Privileged, Peculiar** 897 **Self-Knowledge**

898 Oftentimes, cognitively developed adult human beings have knowledge—or
 899 at least beliefs—about their own attitudes, and they often have further beliefs
 900 about how their attitudes converge with or diverge from their peers. It is
 901 often these higher-order states of mind that motivate agents to reason with
 902 one another in the first place. After all, if one agent believes that there is a
 903 discrepancy between what she believes and what her interlocutor believes,
 904 this can help to explain why she bothers to try and settle the discrepancy
 905 through an interpersonal exchange of reasons.

906 For a hypothetical example, consider two interlocutors: Maya and Roman.
 907 Maya might aim to convince Roman that climate change is an existential
 908 threat to human civilization (note that this is an agent-neutral proposition: I
 909 emphasize the importance of this fact near the end of this section). My claim
 910 now is this: Maya would be in a precarious epistemic position, one that might
 911 undermine the efficiency of her reasoning with Roman, or one that might
 912 even make it better for her to *not* try to reason with Roman about this issue, if
 913 she did not possess **PPSK**.

914 Why so? It is easiest to begin by focusing on *privilege*. Here is the basic
 915 idea: if Maya were not in a systematically superior epistemic position regard-
 916 ing her beliefs about her attitudes than Roman was concerning his beliefs
 917 about Maya's attitudes, then Roman could more easily—i.e., with better epis-
 918 temic grounds—convince Maya that her attitudes already align with his. In
 919 convincing Maya of this, Roman would be providing second-order grounds
 920 for skepticism about Maya's belief that she believes climate change to be an

Evidential unity and cognitive seem conceptually possible, and they might even be achieved by actual agents who are wired to one another's brains in the right sorts of ways. The upshot is that neither party must have "I"-thoughts about their selves and attitudes in the course of interpersonal reasoning nor, for that matter, thoughts about others' selves and attitudes. Instead, by focusing strictly on first-order reasons, they can automatically adjust one another's attitudes. However, Roelofs admits that, for us, "it seems very unlikely...that such a close rapport could persist for very long, or cover very many topics" (2017, 17). We are simply not wired to one another in these ways, at least not with any real consistency. Accordingly, what I say below applies to ordinary agents who lack evidential unification and cognitive vulnerability.

921 existential threat to human civilization. As a result, Maya would not even
 922 bother to reason with Roman about the first-order discrepancy, because her
 923 self-belief would change in such a way that she no longer took there to be
 924 any such discrepancy. Roman might alter Maya's self-belief in good faith by
 925 providing evidence that it is mistaken. But in other cases, Roman might operate
 926 in bad-faith by knowingly supplying Maya with misleading grounds for
 927 the same conclusion. Indeed, if Roman's testimony is a source of evidence all
 928 on its own then, given Maya's lack of privileged access to her own belief, her
 929 epistemic situation upon receiving Roman's testimony is immediately altered
 930 even if Roman supplies no independent evidence in favour of his testimony.
 931 In such cases we could say—perhaps somewhat overdramatically—that Maya
 932 has been taken as Roman's *epistemic hostage*. As an epistemic hostage, Maya
 933 succumbs to Roman's efforts (good faith or otherwise) to convince her that her
 934 self-ascribed attitudes are not really her own. Maya, being falsely convinced
 935 of this, is even cut off from opportunities to reason with agents *other than* Roman
 936 about climate change, given that she has been pre-emptively convinced
 937 that she does not disagree with those—like Roman—who are climate science
 938 deniers.

939 We might construe these situations as threats to Maya's epistemic auton-
 940 omy. I say this because, plausibly, epistemic autonomy is at least partly a
 941 matter of being able to navigate various interpersonal reasoning contexts
 942 without having one's self-conception co-opted too easily by others. Indeed,
 943 this matters even if we are sometimes duped about the *first-order* issues by
 944 clever interlocutors who supply us with misleading evidence at *that* level
 945 of discourse (e.g., misleading statistics suggesting that climate change—of
 946 the anthropogenic variety, at least—is not taking place). An agent who is
 947 convinced by a clever interlocutor that the evidence for climate change is bad
 948 is still an agent who has assessed those reasons for herself and hence has been
 949 mislead on a basis that still deploys her own rational faculties to some degree.
 950 And while it is true that Maya might also deploy her own rational faculties in
 951 assessing Roman's claim that her *self-belief* is wrong, perhaps because Roman
 952 supplied her with good reasons (by *her* lights, at least) to do so, the *result* is
 953 that Maya lacks the self-knowledge that she needs in order to recognize that
 954 there is a discrepancy between her belief about climate change and Roman's
 955 belief about it, and *this* undermines her epistemic autonomy for reasoning
 956 with Roman about climate change itself.

957 Now, as aforementioned, this account of PPSK's instrumental value is
 958 most clearly geared toward *privileged* self-knowledge, since it is an argument

959 about what happens when the epistemic security of Maya's self-beliefs is,
960 as a general matter, no better than that of Roman's perspective unto Maya's
961 mind. But the account can extend to *peculiarity* as well, at least if we construe
962 the relationship between privilege and peculiarity in such a way that Maya's
963 privilege *is due to* the peculiar way in which she knows her own mind (cf.
964 Peterson 2021). For, if her self-beliefs are not generally acquired by a peculiar
965 means that is generally available to her, then nothing prevents individuals
966 like Roman from seizing upon the very same means to acquire knowledge of
967 Maya's mind, and this makes it harder to understand why Maya's self-beliefs
968 are, in general, so epistemically secure that Roman's contrary claims or beliefs
969 do not give Maya strong reason to change what she believes about herself.

970 To bring this point into sharper relief, we can consider a putative foil for my
971 account, namely, Quassim Cassam's *Inferentialist* account of self-knowledge.
972 According to Cassam, both self-knowledge and other-knowledge of agents'
973 attitudes are acquired through inferences. On his view, there remains an
974 epistemic asymmetry between self-knowledge and other-knowledge, but this
975 asymmetry simply "boils down to a difference in the kinds of evidence that are
976 available in the two cases" (Cassam 2014, 150). More specifically, the evidence
977 that one has about one's own attitudes is superior to the evidence that one has
978 about others' attitudes because it includes sensations, memories, and other
979 non-attitudinal mental goings-on that are not so easily accessed by one's peers.
980 Applying this view to interpersonal reasoners like Maya, we might say that
981 Maya's self-knowledge of her attitudes is privileged to some degree even if the
982 same method—inference—is used by both Maya and Roman in coming to
983 form beliefs about Maya's attitudes. So there is nothing peculiar about Maya's
984 route to self-knowledge. But now one might insist that Maya cannot be easily
985 taken as an epistemic hostage by Roman, even though she lacks a peculiar
986 way of knowing herself, simply because she has especially good evidence
987 about herself.

988 However, it could happen that such additional evidence is unavailable to
989 Maya in any number of cases, for what reason can be given for thinking that
990 Maya will always have access to special evidence, given that access to evidence
991 in general is a contingent matter of one's epistemic position relative to a body
992 of information? Peter Carruthers—another prominent Inferentialist—takes
993 it that we have privileged self-knowledge of non-propositional-attitudinal
994 mental states (2011), and contends that this can be used as a basis for inferring
995 our propositional attitudes. However, privileged access to these other mental
996 states can only provide a basis for inferring our propositional attitudes *when*

997 *we are in such mental states*, and yet this itself is a contingent matter. Moreover,
 998 even stipulating that Maya has systematically better evidence about herself
 999 than Roman has about her, we would also need a general assurance that Maya
 1000 infers the correct conclusions from this systematically superior evidence. It
 1001 could happen that Maya has privileged access to the evidence about what she
 1002 herself believes but cannot reliably *utilize* this evidence. At the very least, it
 1003 could happen that she is, in general, no better at utilizing this evidence than
 1004 Roman is at utilizing *his* evidence about Maya's attitudes. Indeed, philosophers
 1005 like Carruthers seem to embrace this point when they claim that Inferentialist
 1006 views best explain failures of self-knowledge.

1007 Finally, Inferentialist views are vulnerable to what I call an *efficiency concern*
 1008 and a *gridlock concern*. The efficiency concern is that, absent peculiar access, it
 1009 could be generally appropriate for Roman to ask Maya to supply the grounds
 1010 for her self-beliefs, and for Maya to ask Roman to do the same, just to be
 1011 sure that they were operating in a case where Maya really did have (and had
 1012 effectively utilized) this superior evidence. Engaging in this second-order
 1013 interpersonal reasoning would significantly slow down their efforts to get to
 1014 the first-order issues, thus rendering interpersonal reasoning about first-order
 1015 issues a less efficient activity. The gridlock concern is that the second-order
 1016 issue might not get resolved at all whenever both parties fail to reach a verdict
 1017 about what Maya believes. One might attempt to circumvent these concerns
 1018 by arguing that Maya's inferences are subpersonal or non-conscious, such that
 1019 she cannot be expected to articulate them to Roman. But inferences that are
 1020 not available for peer-review are also inferences that Maya might be required
 1021 to lower her trust in, thus calling her self-beliefs into question all over again.
 1022 To be sure, if some sort of Inferentialism is true, it may follow that agents like
 1023 Maya often have better evidence and draw better inferences about their own
 1024 attitudes than their interlocutors can draw about her attitudes, but the points
 1025 I have been making suggest that such access will be *worse* for Maya than any
 1026 form of access that renders the special epistemic security of her self-beliefs a
 1027 non-contingent matter.

1028 Now, even though I have been critiquing an Inferentialist rejection of pecu-
 1029 liar access, I want to reiterate a general lesson for all would-be skeptics about
 1030 such access. The lesson is that, if the same method—whether inferential or
 1031 otherwise—is used for acquiring both self-knowledge and other-knowledge,
 1032 then epistemic privilege will seem to be highly contingent. For, if two agents
 1033 can come to know one agent's mind by the same means, then there need be
 1034 no systematic barrier to their doing with equal epistemic pedigree. In the

1035 context of my account, this would mean that there is no strong assurance
1036 that agents are systematically warranted in retaining their self-beliefs when
1037 challenged by their interlocutors. And this, in turn, would mean that there
1038 is no general assurance that interpersonal reasoning *about the world*, rather
1039 than about the interlocutors' minds, can proceed smoothly. The efficiency and
1040 gridlock concerns also generalize: if Maya and Roman share the same method
1041 for arriving at a view about Maya's mind, then Roman might endeavour to
1042 interrogate Maya about whether their current context is one in which she
1043 has exercised the method more effectively, whether the method is inferential
1044 or not. This would slow down and (potentially) gridlock the discourse at the
1045 second-order level. Crucially, though, I am not claiming that the systematic
1046 protection provided by PPSK against these concerns is universal in scope.
1047 For my purposes, PPSK's instrumental value will have been demonstrated if
1048 it is our standard sort of self-knowledge. This would ensure that one is not
1049 *systematically, generally, or universally* vulnerable to innocent-yet-erroneous
1050 self-belief change, bad-faith epistemic hostage-takers, or to the efficiency
1051 and gridlock concerns, thus improving interpersonal reasoning's reliability
1052 and efficiency as a tool in our social-epistemic toolkit for understanding our
1053 shared world.

1054 So goes my account. If correct, it shows that PPSK is instrumentally valuable
1055 for interpersonal reasoning, at least among those who are in a position to
1056 form beliefs about their own attitudes in the first place (again, a child who
1057 has yet to acquire the concept of belief cannot be erroneously convinced that
1058 she *shares a belief* with someone else). Notably, the account applies whether
1059 we imagine interpersonal reasoners as aiming to debate an agent-neutral
1060 proposition or an agent-specific one. I initially described Maya as aiming to
1061 convince Roman that climate change is an existential threat—this being an
1062 agent-neutral proposition—whereupon Roman steers the discourse to the
1063 second-order level in order to convince Maya that she does not really believe
1064 this in the first place. But the content of the proposition was incidental to the
1065 example. Had the proposition's content been agent-specific, e.g., about Maya's
1066 particular climate-focused actions or the belief-desire pairs that rationalize
1067 her actions, Roman might have proceeded in the same way. So, my account has
1068 a broader scope than Siewert's: it applies to agent-specific *and* agent-neutral
1069 interpersonal reasoning.

1073 Addressing Objections

1071 In this section I reinforce my account by addressing four objections. The
 1072 first objection is that legitimate challenges to our self-knowledge are in fact
 1073 quite frequent, and that this provides evidence against the claim that PPSK
 1074 frequently serves as an epistemic shield against erroneous self-belief change
 1075 in our actual interpersonal reasoning practices. The second objection is that
 1076 PPSK does not suffice to ensure that interpersonal reasoning is a reliable
 1077 route for rational attitude adjustments.⁸ The two final objections are specific
 1078 defenses of the claim that factors beyond PPSK can protect interpersonal
 1079 reasoners against erroneous self-belief change in interpersonal reasoning
 1080 contexts.

1081 The first objection turns on familiar cases of self-deception. Self-deception is
 1082 ordinarily taken as a failure of self-knowledge in which an agent self-ascribes
 1083 an attitude that she in fact lacks. Those who take us to have privileged self-
 1084 knowledge surely ought to say something about this familiar phenomenon.
 1085 If one does not take privileged access to be universal in scope, then it is at
 1086 least logically possible to accommodate such cases. Alternatively, one might
 1087 deny the ordinary view of self-deception by arguing that it does not involve
 1088 false self-ascriptions (Bilgrami 2006; Coliva 2016). The apparent trouble for
 1089 my account, however, is that accusations of self-deception are frequent and
 1090 potentially epistemically legitimate in many cases, and yet these might be
 1091 precisely the moves that our interlocutors use in order to convince us that
 1092 our self-beliefs are false. If accusations of self-deception are epistemically
 1093 legitimate and widespread, and if these accusations can spur agents to adjust
 1094 their self-beliefs, then what protection does PPSK really provide here?

1095 To begin my response, I want to reiterate a point from my introductory
 1096 remarks about the dialectical ambitions of this paper, namely, that readers
 1097 need not be convinced that I have unearthed PPSK's actual functional role
 1098 for interpersonal reasoners at this world. Secondly, when I say that PPSK
 1099 provides an epistemic shield against challenges to one's self-beliefs in inter-
 1100 personal reasoning, I do not deny that people might sometimes fail to take
 1101 advantage of this shield—PPSK offers *epistemic* protection that may not be
 1102 *psychologically* appreciated. Beyond these somewhat concessionary responses,
 1103 the devil must reside in the details, since any further response depends on
 1104 how we understand the cases at issue. Thus, consider a case in which Maya

8 These first two objections were put to me by Rachel Cooper.

1105 avows a love of comic books and Roman replies: “you do not love comic books;
 1106 you’ve just tricked yourself into thinking that loving comic books makes you
 1107 interestingly different.”⁹ What might bring Maya to accept this accusation?
 1108 Well, Maya might fixate on the thought that her interlocutor has better evi-
 1109 dence about her mind than Roman has about it. If she wondered about her
 1110 own evidence, and wondered about its inferential role in supporting her self-
 1111 beliefs, she would be supposing her own self-ascription to be vulnerable to
 1112 the same epistemic standards that Roman uses to evaluate her self-beliefs. If
 1113 her self-knowledge is peculiar, however, she will not fixate on this thought,
 1114 because her self-ascription is *not* based on the same epistemic standards.

1115 In fact, our actual manner of proceeding tends to bear this out: one’s inter-
 1116 locutor judges one to be self-deceived about one’s love of comic books, and
 1117 one responds *not* by attempting (and possibly failing) to offer higher-quality
 1118 evidence *about what one believes*, but by offering reasons about *why comic*
 1119 *books are loveable*. Indeed, one possible *explanation* of privilege and peculiari-
 1120 ty is that one’s own take on the reasons for or against adopting some attitude
 1121 (typically) determine one’s adoption of it. And if one self-ascribes this attitude
 1122 with full knowledge of the first-order reasons that one takes to support it,
 1123 one is entitled to make this self-ascription even if other people have evidence
 1124 contravening one’s self-ascription (Bilgrami 2006; Coliva 2016).¹⁰

1125 Moreover, if we have PPSK, *other* challenges to our self-beliefs may also be
 1126 illegitimated, these being challenges where other agents do not accuse us of
 1127 being self-deceived but, rather, as having made innocent (or “brute”¹¹) errors
 1128 about ourselves—errors that could only be made on the basis of innocent
 1129 inferential or observational mistakes.

1130 Here is another, final sense in which the devil is in the details. The objection
 1131 under consideration is that accusations of self-deception are common, and
 1132 that these accusations might frequently lead to (reasonable) changes in one’s
 1133 self-beliefs. However, while such cases may indeed be common, they may
 1134 only be common in the sense that *all of us* are *occasionally* susceptible to
 1135 them. On this explanation of their commonality, no single agent is liable to be
 1136 the reasonable target of an overwhelmingly large number of self-deception
 1137 accusations. There is something suspicious about anyone, even one’s therapist,

9 I owe this example to Rachel Cooper.

10 Compare Schwengerer’s verdict on two cases he discusses (2021, 77). What I may owe my interlocutor, in this case, is an explanation of how my actions fail to live up to my self-ascribed attitude, *not* an explanation to the effect that the evidence shows that I have this attitude.

11 For the operative notion of brute error, see Burge (1996) and Bar-On (2004).

1138 who would unrelentingly accuse one of self-deception across myriad cases
 1139 by saying things like “you do not believe that-P, nor hope that-Q, nor desire
 1140 to φ , nor love S....” This suspicion may well reflect a fact about us: that we
 1141 have enough PPSK to be reasonable in *not* giving in to too many accusations
 1142 of self-deception—accusations which, if legitimate, would force us to change
 1143 our self-beliefs.

1144 The second objection to my account is that PPSK does not improve the
 1145 reliability of interpersonal reasoning even if it provides us with epistemic
 1146 warrant to ignore (many) accusations of mistaken self-belief. Cases in favour
 1147 of this objection are easy enough to set up. For example, maybe Maya claims
 1148 that climate change is an existential threat to human civilization and Roman
 1149 gives insufficient epistemic uptake to her assertion because he is prejudiced
 1150 against women. Indeed, in this case, Maya may be the victim of a “testimonial
 1151 injustice” (Fricker 2007). But I want to offer two observations. First, although
 1152 the factors preventing Roman from reasoning with Maya have nothing to do
 1153 with Maya’s self-knowledge or Roman’s beliefs about Maya’s self-perspective,
 1154 this does not change the fact that Maya would have an *additional* problem
 1155 on her hands if Roman were generally in an epistemic position to make
 1156 Maya erroneously change her self-beliefs. Second, to the extent that Roman’s
 1157 prejudiced behaviour does not prevent Maya from knowing herself, she is still
 1158 in a position to congregate with less prejudiced individuals and to reason with
 1159 *them* (or even to reason with Roman *indirectly* by reasoning with someone
 1160 that Roman is *not* prejudiced against, and getting that individual to convey
 1161 Maya’s reasons to Roman). This point also applies to another concern, namely,
 1162 that Roman might simply say that *he* agrees with Maya when *he* does not (this
 1163 being an inverse version of the epistemic hostage-taking tactic). Maya may
 1164 not be able to rationally challenge this claim if Roman has PPSK, unless she
 1165 has reason to deem him insincere, since she will then have strong reason to
 1166 take Roman at his word. Once again, though, this would not put Maya in the
 1167 position of being made to form a false belief about what she herself believes
 1168 about the world, and hence she would not be prevented from discoursing
 1169 with other agents about the contents of her beliefs about the world.

1170 I now address two objections to the effect that something other than PPSK
 1171 can explain why we are protected against epistemic hostage-taking. According
 1172 to the first objection, what *really* protects Maya against Roman’s nefarious
 1173 machinations is the same thing as what explains her *first-person authority*,
 1174 where what explains *this* is something other than PPSK. Roughly, “first-person
 1175 authority” denotes two claims: (1) it is epistemically rational to presume

1176 the truth of speakers' present-tense self-ascriptions of mental states, and (2)
1177 it is typically epistemically irrational to interrogate the epistemic grounds
1178 of speakers' present-tense self-ascriptions (hereafter just 'self-ascriptions').¹²
1179 Now consider an "expressivist" explanation of first-person authority which
1180 contends that speakers' self-ascriptions ought to be presumed true and be
1181 insulated from requests for epistemic support because they express and thus
1182 *show* the self-ascribed mental states to one's hearers (Bar-On 2004). This
1183 explanation is available even if speakers do not *also* possess PPSK of the
1184 mental states that their self-ascriptions express. The objection, then, is that
1185 Maya's first-person authority gives Roman a strong reason not to challenge
1186 most of her self-ascriptions, such that PPSK is explanatorily superfluous
1187 in explaining why Roman is not likely to give Maya an erroneous basis for
1188 changing her self-beliefs.

1189 Now, for all I have said, Maya's self-ascriptions may be first-person authori-
1190 tative in virtue of what they express, whether or not Maya also has PPSK of
1191 what they express. Nevertheless, I argue that without *also* possessing PPSK,
1192 Roman could purposefully *ignore* the first-person authority of Maya's self-
1193 ascriptions in a bid to convince her that her attitudes converge with rather
1194 than diverge from Roman's. He might (rightly) take Maya to have expressed
1195 her first-order belief through a self-ascription but still claim that her self-
1196 belief is false. Hence, PPSK protects Maya against being manipulated by bad
1197 faith interlocutors who ignore her first-person authority, *however* that is to be
1198 explained, because PPSK ensures the general (and systematically superior)
1199 reliability of her self-beliefs relative to Roman's beliefs about her attitudes.
1200 PPSK is what gives Maya an epistemic warrant for holding steadfast against
1201 his machinations, even if he was already unjustified in challenging her self-
1202 ascription challenged her self-ascription in the first place.¹³ Moreover, PPSK
1203 protects Maya against erroneous self-belief change even if Roman, innocently,
1204 fails to recognize that her self-ascription expresses the very attitude that it is
1205 about.

1206 The final objection to my argument is that Maya can get away with merely
1207 *assuming* that she generally has PPSK, such that she is generally entitled to
1208 not defer to interlocutors who challenge her self-beliefs (whether in good
1209 or bad faith). More substantively, it might be argued that Maya possesses a

12 See Doyle (2021) and Winokur (2022) for more precise articulations of these claims.

13 I take expressivism to have brighter prospects than Schwengerer (2021) does, though I also agree with him that not *everything* epistemically interesting about mental state discourse can be explained by first-order phenomena, hence the account given in this paper.

1210 distinctively strong *practical* warrant for holding steadfast when faced with
 1211 accusations of mistaken self-belief, even if she lacks a distinctively strong
 1212 *epistemic* warrant for doing so.

1213 The trouble with this objection is that it is hard to see what could ground
 1214 Maya's practical warrant for holding steadfast if it is not really, at bottom, the
 1215 same as (or itself grounded by) epistemic warrant for doing so. This is because a
 1216 *merely* practical warrant here would go against her epistemic wellbeing in any
 1217 number of cases. Specifically, if she did not systematically know herself better
 1218 than others know her, then she *would* often—perhaps even typically—have
 1219 an epistemic reason to discourse with others about whether her self-beliefs
 1220 are true, and this would be in tension with her practical warrant for avoiding
 1221 such discourse. In other words, it is only if Maya really has PPSK, thus having
 1222 epistemic warrant for holding on to her self-beliefs, that holding steadfast
 1223 against her interlocutors' countervailing assertions does not inadvertently
 1224 prevent her from indulging many epistemically legitimate disagreements
 1225 about what her attitudes are. It is only if she really has PPSK that *not* entering
 1226 these disagreements is by and large good for her to do.

1227 It might now be complained, relatedly, that I have merely established the
 1228 importance of an especially strong epistemic warrant for our self-beliefs, but
 1229 that this need not amount to PPSK. In other words, Maya might have an
 1230 especially strong epistemic warrant for her self-beliefs, but these self-beliefs
 1231 need not be especially reliably *true*.¹⁴ Indeed, such warrant may also suffice
 1232 for avoiding the efficiency and gridlock concerns described in section 2. But
 1233 I think a similar response applies here. For, if Maya has especially strong
 1234 epistemic warrant for her self-beliefs but this warrant does not amount to
 1235 self-knowledge in at least most of the cases in which she possesses this war-
 1236 rant, then in any number of cases she will still miss out on an epistemic
 1237 good—that of *true* warranted self-belief—whenever she declines to engage
 1238 with interlocutors who claim that her self-beliefs are false. Moreover, it is
 1239 hard to understand how she could possess this special epistemic warrant
 1240 for her self-beliefs if she did not actually possess self-knowledge in most of
 1241 those same cases. After all, this would be tantamount to having epistemic
 1242 warrant for self-beliefs that were not correspondingly likely to be true, and

14 Thanks to an anonymous reviewer for this objection.

1243 this systematic mismatch between truth and warrant would call the warrant
1244 itself into question.¹⁵

1244 **4 Interpersonal Reasoning and Epistemic Control**

1246 In sections 2–3 I argued that **PPSK** provides us with widespread (even if not
1247 universal) protection against situations in which others provide epistemic
1248 reasons for us to change our self-beliefs, whether our interlocutors are operat-
1249 ing innocently or in bad faith, and that this helps to ensure the efficiency and
1250 reliability of interpersonal reasoning. In this section I show that my account
1251 complements another recent account of **PPSK**'s instrumental value.

1252 According to Jared Peterson, **PPSK** is instrumentally valuable because it
1253 facilitates “epistemic control,” which is a matter of being able to “keep private
1254 or disclose particular facts about one’s mind to others” (2021, 368). Take
1255 privacy first. If you have **PPSK**, then you can reliably conceal your attitudes
1256 from others. For example, a teacher might fail to motivate a student’s learning
1257 if the student knows that the teacher is pessimistic about the student’s progress.
1258 But if the teacher has **PPSK** of her pessimism, then she has greater epistemic
1259 control over whether the student discovers this. Therefore, the teacher has
1260 greater control over the student’s motivation to continue studying. For an
1261 example about disclosing rather than concealing one’s mental states, Peterson
1262 says that “[a]n estranged lover might want a former partner to know in a
1263 highly epistemically secure manner that she still loves him” (2021, 369).

1264 He also says that epistemic control:

1265 [...] allows societies to function in a much more productive, orga-
1266 nized, and amicable way. When we accomplish group objectives
1267 in an efficient and peaceful manner we do so in large part by
1268 keeping private that which would be counterproductive to the
1269 group’s efforts, and/or revealing our thoughts, beliefs, desires, etc.
1270 that are valuable for other members of a group to know. (Peterson
1271 2021, 371)

1272 Peterson and I are both happy to emphasize the social importance of **PPSK**.
1273 I am also willing to say that **PPSK** provides a way to reliably disclose one’s

15 This response is similar in structure to one pursued by Davidson (1991) regarding perceptual belief warrant, though I believe that the strength of our warrants for perceptual beliefs and self-beliefs differ.

attitudes to others. However, Peterson does not also acknowledge the additional possibility of, e.g., *expressing one's love itself* as a reliable way of putting one's former partner in a secure position with respect to one's mind, where this expressive capacity may or may not depend on an agent's self-knowledge.

More significantly, though, I submit that PPSK's role as a shield against erroneous self-belief change is independent of its role in enabling one to disclose or conceal one's attitudes from others. To be able to better conceal one's attitudes is to prevent others from discovering what attitudes one has, but this may not matter to interlocutors who do not care (or are simply mistaken) about the facts and, instead, aim to convince you of a certain belief about yourself. Similarly, having an especially epistemic secure way of disclosing your attitudes is something that interlocutors could ignore (as argued in section 3 when discussing first-person authority). Thus, one may be tempted to deny that the instrumental value of PPSK for interpersonal reasoning is a matter of epistemic control.

However, one might be just as easily inclined to regard this as an instance of epistemic control after all, since my account claims that agents with systematically superior knowledge of their self-beliefs thereby exercise greater control over their social-epistemic lives as interpersonal reasoners. Accordingly, the lesson to be drawn may be that we ought to broaden our view of PPSK's contribution to epistemic control, such that epistemic control encompasses (i) control over which attitudes one discloses to others,¹⁶ (ii) control over which attitudes one conceals from others, and (iii) control over which attitudes one is able to self-attribute, with especially secure epistemic warrant, in the face of disagreement about one's attitudes, while attempting to reason with others.

5 Conclusion

I have argued that privileged and peculiar self-knowledge contributes to our capacity for interpersonal reasoning about the world around us. To the extent that agents possess PPSK of their attitudes, interpersonal reasoning is a more reliable route to discursively navigating our shared world, and this explains one way in which PPSK is instrumentally valuable.

For the record, I happen to believe that phenomena like epistemic hostage taking are not widespread, and I regard PPSK as at least a partial explanation of this fact. I take myself, therefore, to have contributed to the debate

¹⁶ Again, if this particular capacity requires PPSK at all.

1308 between PPSK-skeptics and PPSK-realists not *merely* by illuminating the de-
 1309 bate's stakes, but *also* by taking an anti-skeptical stand within that debate.
 1310 This being said, I reiterate that others may not be persuaded to go as far as me
 1311 in this regard, such that the core contribution of this paper is best viewed as
 1312 an account of how being a PPSK-skeptic or PPSK-realist should affect one's
 1313 corresponding conception of our interpersonal reasoning competencies.*

1314 Benjamin Winokur

1315  0000-0002-0845-9460

1316 University of Macau

1317 bwinokur@um.edu.mo

1318 References

- 1319 ANDREOTTA, Adam J. 2021. "Confabulation Does Not Undermine Introspection for
 1320 Propositional Attitudes." *Synthese* 198(5): 4851–4872, doi:10.1007/s11229-019-
 1321 02373-9.
- 1322 —. 2022. "More than Just a Passing Cognitive Show: A Defence of Agentialism About
 1323 Self-Knowledge." *Acta Analytica* 37(3): 353–373, doi:10.1007/s12136-021-00492-y.
- 1324 BAR-ON, Dorit. 2004. *Speaking My Mind: Expression and Self-Knowledge*. Oxford:
 1325 Oxford University Press, doi:10.1093/0199276285.001.0001.
- 1326 BILGRAMI, Akeel. 2006. *Self-Knowledge and Resentment*. Cambridge, Massachusetts:
 1327 Harvard University Press, doi:10.2307/j.ctv1nznfgcn.
- 1328 BOYLE, Matthew. 2011. "Transparent Self-Knowledge." *Proceedings of the Aristotelian*
 1329 *Society, Supplementary Volume* 85: 223–241, doi:10.1111/j.1467-8349.2011.00204.x.
- 1330 —. 2019. "Transparency and Reflection." *Canadian Journal of Philosophy* 49(7):
 1331 1012–1039, doi:10.1080/00455091.2019.1565621.
- 1332 BRANDON, Robert B. 1994. *Making It Explicit: Reasoning, Representing, and Discursive*
 1333 *Commitment*. Cambridge, Massachusetts: Harvard University Press.
- 1334 BURGE, Tyler. 1996. "Our Entitlement to Self-Knowledge." *Proceedings of the Aris-*
 1335 *totelian Society* 96: 91–116. Reprinted in Burge (2013, 68–87), doi:10.1093/aristote-
 1336 lian/96.1.91.

* Thanks are due to Adam Andreotta for his comments on an earlier version of this paper, and to Rachel Cooper for reading and commenting on a shortened version of it at the 2021 American Philosophical Association Pacific Meeting. Thanks are also due to the other audience members at that meeting, to audience members at the 2021 Canadian Philosophical Association Meeting, and to audience members at the 2022 Meeting of the International Social Ontology Society. Finally, special thanks are due to three anonymous reviewers at *Dialectica* for their many useful comments.

- 1337 —. 2013. *Cognition Through Understanding: Self-Knowledge, Interlocution, Reasoning,*
 1338 *Reflection. Philosophical Essays, Volume 3.* Philosophical Essays n. 3. Oxford: Oxford
 1339 University Press, doi:10.1093/acprof:oso/9780199672028.001.0001.
- 1340 BYRNE, Alex. 2018. *Transparency and Self-Knowledge.* Oxford: Oxford University Press,
 1341 doi:10.1093/oso/9780198821618.001.0001.
- 1342 CAMPBELL, Lucy. 2018. “Self-Knowledge, Belief, Ability (and Agency?).” *Philosophical*
 1343 *Explorations* 21(3): 333–349, doi:10.1080/13869795.2018.1426779.
- 1344 CARRUTHERS, Peter. 2011. *The Opacity of Mind: An Integrative*
 1345 *Theory of Self-Knowledge.* Oxford: Oxford University Press,
 1346 doi:10.1093/acprof:oso/9780199596195.001.0001.
- 1347 CASSAM, Quassim. 2014. *Self-Knowledge for Humans.* Oxford: Oxford University Press,
 1348 doi:10.1093/acprof:oso/9780199657575.001.0001.
- 1349 COLIVA, Annalisa. 2016. *The Varieties of Self-Knowledge.* Innovations in Philosophy.
 1350 London: Palgrave Macmillan, doi:10.1057/978-1-137-32613-3.
- 1351 DAVIDSON, Donald. 1991. “Epistemology Externalized.” *Dialectica* 45(2/3): 191–202.
 1352 Reprinted in Davidson (2001), doi:10.1111/j.1746-8361.1991.tb00986.x.
- 1353 —. 2001. *Subjective, Intersubjective, Objective: Philosophical Essays, Volume 3.* Oxford:
 1354 Oxford University Press, doi:10.1093/0198237537.001.0001.
- 1355 DOYLE, Casey. 2021. “There’s Something About Authority.” *Journal of Philosophical*
 1356 *Research* 46: 363–374, doi:10.5840/jpr2021816169.
- 1357 FRICKER, Miranda. 2007. *Epistemic Injustice: Power and the Ethics of Knowing.* Oxford:
 1358 Oxford University Press, doi:10.1093/acprof:oso/9780198237907.001.0001.
- 1359 GOPNIK, Alison. 1993. “How We Know Our Minds: The Illusions of First-Person
 1360 Knowledge of Intentionality.” *Behavioral and Brain Sciences* 16(1): 1–14, doi:10.1
 1361 017/S0140525X00028636.
- 1362 GRECO, John. 2019. “The Social Value of Reflection.” in *Thinking about Oneself: The*
 1363 *Place and Value of Reflection in Philosophy and Psychology*, edited by Waldomiro
 1364 J. SILVA-FILHO and Luca TATEO, pp. 45–57. Philosophical Studies Series n. 141.
 1365 Cham: Springer, doi:10.1007/978-3-030-18266-3_4.
- 1366 KEELING, Sophie. 2018. “Confabulation and Rational Obligations for Self-Knowledge.”
 1367 *Philosophical Psychology* 31(8): 1215–1238, doi:10.1080/09515089.2018.1484086.
- 1368 —. 2019a. “The Transparency Method and Knowing Our Reasons.” *Analysis* 79(4):
 1369 613–621, doi:10.1093/analys/anz031.
- 1370 —. 2019b. “Knowing Our Reasons: Distinctive Self-Knowledge of Why We Hold Our
 1371 Attitudes and Perform Actions.” *Philosophy and Phenomenological Research* 102(2):
 1372 318–341, doi:10.1111/phpr.12655.
- 1373 KOREŇ, Ladislav. 2023. “Have Mercier and Sperber Untied the Knot of Human Rea-
 1374 soning?” *Inquiry* 66(5): 849–862, doi:10.1080/0020174X.2019.1684988.
- 1375 LEVINE, Steven M. 2009. “Expressivism and I-Beliefs in Brandom’s *Making it Explicit.*”
 1376 *International Journal of Philosophical Studies* 17(1): 95–114, doi:10.1080/096725
 1377 50802614786.

- 1378 MARCUS, Eric and SCHWENKLER, John. 2019. "Assertion and Transparent
1379 Self-Knowledge." *Canadian Journal of Philosophy* 49(7): 873–889,
1380 doi:10.1080/00455091.2018.1519771.
- 1381 MÜLLER, Andreas. 2019. "Reasoning and Normative Beliefs: Not Too Sophisticated."
1382 *Philosophical Explorations* 22(1): 2–15, doi:10.1080/13869795.2018.1523449.
- 1383 MUSHOLT, Kristina. 2015. *Thinking about Oneself: From Nonconceptual Content to the*
1384 *Concept of a Self*. Cambridge, Massachusetts: The MIT Press, doi:10.7551/mitpre
1385 ss/9780262029209.001.0001.
- 1386 NGUYEN, A. Minh. 2015. "What Good is Self-Knowledge?" *Journal of Philosophical*
1387 *Research* 40: 137–154, doi:10.5840/jpr2015111656.
- 1388 PARENT, Ted. 2017. *Self-Reflection for the Opaque Mind: An Essay in Neo-Sellarsian*
1389 *Philosophy*. London: Routledge, doi:10.4324/9781315618449.
- 1390 PETERSON, Jared. 2021. "The Value of Privileged Access." *European Journal of Philoso-*
1391 *phy* 29(2): 365–378, doi:10.1111/ejop.12594.
- 1392 ROELOFS, Luke. 2017. "Rational Agency without Self-Knowledge: Could 'We' Replace
1393 'I'?" *Dialectica* 71(1): 3–33, doi:10.1111/1746-8361.12169.
- 1394 SCHWENGERER, Lukas. 2021. "Beliefs Over Avowals: Setting Up the Discourse on
1395 Self-Knowledge." *Episteme* 18(1): 66–81, doi:10.1017/epi.2018.56.
- 1396 SHOEMAKER, Sydney S. 1988. "On Knowing One's Own Mind." in *Philosophical Per-*
1397 *spectives 2: Epistemology*, edited by James E. TOMBERLIN, pp. 183–209. Atascadero,
1398 California: Ridgeview Publishing Co. Reprinted in Shoemaker (1996, 25–49),
1399 doi:10.2307/2214074.
- 1400 —. 1996. *The First-Person Perspective and Other Essays*. Cambridge: Cambridge Uni-
1401 versity Press, doi:10.1017/cbo9780511624674.
- 1402 SIEWERT, Charles. 2003. "Self-Knowledge and Rationality: Shoemaker on
1403 Self-Blindness." in *Privileged Access: Philosophical Accounts of Self-Knowledge*,
1404 edited by Brie GERTLER. Ashgate Epistemology and Mind Series. London:
1405 Routledge, doi:10.4324/9781315245997.
- 1406 SORGIOVANNI, Ben. 2019. "The Agential Point of View." *Pacific Philosophical Quarterly*
1407 100(2): 549–572, doi:10.1111/papq.12263.
- 1408 STRIJBOS, Derek W. and DE BRUIN, Leon C. 2012. "Making Folk Psychology Explicit:
1409 The Relevance of Robert Brandom's Philosophy for the Debate on Social Cognition."
1410 *Philosophia* 40(1): 139–163, doi:10.1007/s11406-010-9288-z.
- 1411 VALARIS, Markos. 2018. "Self-Knowledge." in *The Philosophy of Knowledge: A History,*
1412 *Volume 4: Knowledge in Contemporary Philosophy*, edited by Stephen Cade HET-
1413 HERINGTON and Markos VALARIS, pp. 155–174. London: Bloomsbury Academic,
1414 doi:10.5040/9781474258814.ch-008.
- 1415 WINOKUR, Benjamin. 2021a. "Critical Reasoning and the Inferential Transparency
1416 Method." *Res Philosophica* 98(1): 23–42, doi:10.11612/resphil.1967.
- 1417 —. 2021b. "Inference and Self-Knowledge." *Logos & Episteme* 12(1): 77–98, doi:10.584
1418 0/logos-episteme20211214.

- 1419 —. 2022. “There is Something to the Authority Thesis.” *Journal of Philosophical Re-*
1420 *search* 47: 115–132, doi:[10.5840/jpr202298189](https://doi.org/10.5840/jpr202298189).

PROOF

1421 Avner Baz’s Ordinary Language
1422 Challenge to the Philosophical
1423 Method of Cases

PAUL O. IRIKEFE

1424 Avner Baz argues that the philosophical method of cases presupposes a
1425 problematic view of language and linguistic competence, namely, what
1426 he calls “the atomistic-compositional view”. Combining key elements of
1427 social pragmatism and contextualism, Baz presents a view of language
1428 and linguistic competence, which he takes to be more sensitive to the
1429 open-endedness of human language. On this view, there are conditions
1430 for the “normal” and “felicitous” use of human words, conditions that
1431 Baz thinks are lacking in the context of the philosophical method of
1432 cases, and which make the question that philosophers are prone to ask
1433 in that context and the answers they give to that question to be pointless.
1434 However, in this paper, I argue as follows. First, Baz’s conditions for
1435 the “normal” and “felicitous” use of human words are in tension with
1436 the open-endedness of human language and the use of human words.
1437 Second, it is not even clear that those conditions are really missing in the
1438 context of the philosophical method of cases. And third, even if we grant
1439 that those conditions are missing in that context, this does not licence his
1440 damning conclusion on the philosophical method of cases since we are
1441 not forced to embrace the view of language and linguistic competence
1442 on which that damning conclusion is plausible. This last move is secured
1443 by advancing and defending a skill or virtue-based view of language and
1444 linguistic competence inspired by the later work of Donald Davidson.

1445 The philosophical method of cases (henceforth, PMOC) arguably plays some
1446 role in how philosophers investigate issues of great philosophical interest like
1447 knowledge, free will, and reference.¹ In this practice, a philosopher would

1 There is evidence of the use of the method outside the Western tradition of philosophy (Boh 1985).

1448 describe a certain scenario, whether real or hypothetical, and invite us to
 1449 say whether the case so described would count as falling under the relevant
 1450 property or term or concept under investigation. The judgement formed on
 1451 the described scenario is then enlisted in arguing for or against certain philo-
 1452 sophical views.²

1453 The question then is, what linguistic competence guides this practice? In
 1454 some very illuminating works, Avner Baz (2016, 2017) argues that the PMOC
 1455 presupposes a problematic view of language and linguistic competence, what
 1456 he calls the atomistic-compositional view. The atomistic-compositional view
 1457 as he presents it is presupposed by defenders of the method in mainstream
 1458 analytic philosophy and critics of the method, including experimental philoso-
 1459 phers. Combining key elements of social pragmatism and contextualism, Baz
 1460 presents what we might call a “social pragmatic view of language”, a view
 1461 he thinks enjoys better empirical support and is more sensitive to the open-
 1462 endedness of human language. On this view, there are “normal” and “felic-
 1463 itous” conditions for the use of words and human language, conditions he
 1464 takes to be lacking in the context of the PMOC and the questions philosophers
 1465 are prone to ask in that context such as: “Does *X* know *Y*?”

1466 However, in this paper, I argue as follows. First, Baz’s conditions for the
 1467 normal and felicitous use of words and language stand in tension with the
 1468 open-endedness of words and language. Second, it is not even clear that those
 1469 conditions are really missing in the context of the PMOC. And third, even if
 1470 we grant that those conditions are missing in that context, this does not licence
 1471 any damning conclusion on the PMOC since we are not forced to embrace the
 1472 view of language and linguistic competence on which that conclusion seems
 1473 plausible. This last move is secured by advancing and defending a skill-based
 1474 view of language and linguistic competence inspired by Donald Davidson
 1475 (1986).

1476 The paper proceeds as follows. In section 1, I discuss what Baz calls the
 1477 “minimal assumption” about language which he says is presupposed by both
 1478 armchair philosophers and their experimental counterparts. I show that the
 1479 assumption expresses two worries. The first is the correctness worry and the
 1480 second is about the kind of linguistic competence we rely on in the PMOC,

2 In recent times, there have been serious discussions about the evidential status of these judgements, in particular, whether this status is due to their being intuitive (Cappelen 2012, 2014; Deutsch 2015; Earlenbaugh and Molyneux 2009; Ichikawa and Jarvis 2009; Irikefe 2020; Williamson 2007). I would set aside this issue in this paper by staying neutral about the evidential nature of these judgements since nothing here hangs on it.

1481 which he calls the “atomistic-compositional” view. I briefly respond to the
1482 first worry, and I indicate that the second worry is more pressing and would
1483 therefore be of present concern. In section 2, I discuss Baz’s social pragmatic
1484 view of language and linguistic competence, which he takes to have better
1485 empirical support than the atomistic compositional view. I explore some of
1486 the ingredients of the social pragmatic view, its negative implications for
1487 the PMOC and why we might worry that some aspects of the view do not
1488 seem consistent with recognisable features of the PMOC and the nature of
1489 human language itself. In section 3, I explore how we might look to defend or
1490 rely on the PMOC without any problematic assumptions about language and
1491 linguistic competence and without either the atomistic compositional view
1492 or Baz’s social pragmatic view. I end the paper by showing how the present
1493 defence of the PMOC meshes with a broader trend in the epistemology of
1494 philosophy and lends independent support to it.

1495 **1 The Atomistic-Compositional View of Language and the** 1496 **Philosophical Method of Cases**

1497 The philosophical method of cases is a standard practice in analytic philosophy.
1498 A philosopher wants to argue for or against certain views about knowledge,
1499 causation, free will or moral permissibility. An imaginary scenario is described,
1500 and we are asked whether or not a certain property, term or concept obtains
1501 in the described scenario. For example, in Gettier’s 10-coin case, we are asked
1502 the question whether the protagonist in the described scenario knows some
1503 particular proposition, that is, whether the protagonist knows that the man
1504 who will get the job has ten coins in his pocket ([Gettier 1963](#)).

1505 According to Baz ([2016](#), [2017](#)), the method depends on a “minimal as-
1506 sumption” about language to get off the ground, namely, the assumption that
1507 questions like that as presented in the context of the PMOC are

1508 in principle, in order—in the simple sense that they are clear
1509 enough and may be answered correctly or incorrectly—and that,
1510 as competent speakers, we ought to understand those questions
1511 and be able to answer them correctly, just on the basis of the
1512 descriptions of the cases and our mastery of the words in which
1513 the questions are couched. ([Baz 2017](#), 6)

1514 We can distinguish two kinds of worries in the minimal assumption. The
1515 first one is the correctness worry, namely, the worry whether the questions
1516 at stake in the method of cases can be answered correctly or incorrectly,
1517 rightly or wrongly, and what the ontological status of such answers might be
1518 like, precisely whether these answers would be about concepts or the world
1519 independent of concepts (Baz 2017, 6). Baz links this worry with what he calls
1520 the “representational-referential” view of language and traces it to Timothy
1521 Williamson (2007), Herman Cappelen (2012) and Frank Jackson (2011). On
1522 this view, *the primary function* of language at any given moment or as he puts
1523 it “the fundamental aim of (all?!) discourses” (Baz 2017, 74, n.6) is to say true
1524 or false things about the world. Although this is not the worry I intend to
1525 address in this paper, I believe that friends of the PMOC do not need to commit
1526 themselves to any problematic assumption here. On the contrary, I think *pace*
1527 Baz, what they need to hold is that *among other things that human language*
1528 *is for*, human language is used to say true or false things about the world (I
1529 would return to this in section 3). In the same vein, friends of the method may
1530 not need to settle the issue of what the answers to the questions at stake in the
1531 method of cases would be true of, whether they would be true of our concepts
1532 or items in the world existing independently of our conception of them. As
1533 Ernest Sosa noted: “We can conduct our controversies, for example, just in
1534 terms of where the truth lies with regard to them, leaving aside questions of
1535 objectual ontology” (Sosa 2007, 100–101).

1536 The second worry in the minimal assumption is the more pressing one.
1537 And it is the one I wish to address in this paper. It says that

1538 as competent speakers, we ought to understand those questions
1539 [i.e., the questions at stake in the method of cases] and be able to
1540 answer them correctly, just on the basis of the descriptions of the
1541 cases and our mastery of the words in which the questions are
1542 couched. (Baz 2017, 6)

1543 Baz notes that this is an assumption about language that derives from and
1544 is dependent on the atomistic-compositional view of language. In this view,
1545 the meaning of the whole of an utterance comes from the fixed meaning of
1546 the parts of that utterance. Baz traces the atomistic-compositional view of
1547 language to Jackson (2011), who presents it as the linguistic competence that
1548 the method depends on. Jackson says that how a sentence like “it is raining
1549 outside” represents things is a

1550 function of the representational contents of its parts and how
1551 they are combined.³ Moreover, we have a grasp of the representa-
1552 tional contents of these parts, and of the way various modes of
1553 combination into sentences generate representational structures
1554 whose contents are a function of the contents of their parts and
1555 the way the parts are put together. (Jackson 2011, 472)⁴

1556 In Jackson (2011), this view of language and linguistic competence goes side
1557 by side with a view of conceptual competence. On this view, in learning
1558 philosophically significant terms like “knowledge” we are latching onto the
1559 pattern or rule or categorisation of “knowledge.” Thus, he says:

1560 How did we acquire the word “knowledge”? We came across lots
1561 of examples. We were told a bit about what mattered. Perhaps, we
1562 were simply instructed that if it is false, it cannot be knowledge.
1563 At some point we latched onto *the* pattern. (Jackson 2011, 474)

1564 This rule or pattern on Jackson's view in turn guides our knowledge ascrip-
1565 tions, that is, it enables us to say whether or not the protagonist in a Gettier
1566 text knows or does not know a given proposition.

1567 In the next section, I consider Baz's argument that the atomistic-
1568 compositional view of language is problematic and his argument that in
1569 the context of the PMOC the conditions for the normal and felicitous use of
1570 words and language are lacking. As we shall see too, Baz takes himself to
1571 be establishing a demarcation of the boundary of linguistic sense, one that
1572 makes clear that the PMOC is outside that boundary and that the questions
1573 philosophers are prone to ask in that context are fundamentally problematic.

3 Compare the atomistic-compositional view with the view of Paul Elbourne: “Suppose you are interpreting an uttered sentence. In a series of extremely intricate processes that are largely subconscious, you access the sentence's words in your mental lexicon and find their meanings; you work out the intended sense of any ambiguous words it might contain; you work out the references of indexicals in the sentence; you work out the sentence's syntactic structure and resolve any structural ambiguities there may be; and you combine the contents of the words in the compositional semantics... If implicit content is not mediated by means of covert indexicals (and thus covered by the second step mentioned above), you add some of this too. Finally, you have worked out the content of the sentence, as uttered on that occasion” (Elbourne 2011, 131; cited in Baz 2017).

4 Although differently expressed, Baz identifies Williamson as holding this view as well: “[E]xpressions refer to items in the mostly non-linguistic world, *the reference of complex expressions is a function of the reference of its constituents*, and the reference of a sentence determines its truth value” (Williamson 2007, 281, emphasis mine).

1572 **Baz's Social Pragmatic View of Language and Linguistic**
 1575 **Competence**

1576 The way Baz shows the atomistic-compositional view of language to be prob-
 1577 lematic is by presenting and defending an alternative view of language that
 1578 he takes to enjoy better empirical support. He finds support from a scientific
 1579 study of how children acquire their first natural language (Bartsch and Well-
 1580 man 1995). But Karen Bartsch and Henry Wellman were not interested in
 1581 natural language acquisition for its own sake. More specifically, they were
 1582 tracking the natural development in the use of belief-desire terms in children
 1583 between ages one and a half to six years. Six of these children are boys and four
 1584 are girls. One of them is African American and the others are not. Because of
 1585 their interest, Bartsch and Wellman were necessarily selective. They were cod-
 1586 ing only for terms expressing genuine psychological reference, where this is
 1587 judged so if with respect to a suitable context it referred to psychological states
 1588 like desire, belief, or knowledge. As a result, they discounted conversational
 1589 use of belief-desire terms like when a child says “you know what?” when
 1590 seeking to get someone’s attention; repetition of phrases uttered by someone
 1591 else, for example, a mother saying “tell him you know where it is,” to which
 1592 the child responds “I know where it is”, and so on.

1593 For present purposes, let us focus on what the study uncovered about
 1594 the term “knows” and its cognates. The authors found (as Baz pointed out)
 1595 that the word “knows” and its cognates do not admit of a simple formula.
 1596 More specifically, they found that children use “knows” and its cognates to
 1597 refer to instances of belief “felt to be justified, assumed to be true, or that
 1598 enjoys markedly higher conviction than one described by *think*” (Bartsch
 1599 and Wellman 1995, 40). Later on in their development, they use it to refer to
 1600 “situations involving successful actions or to correct statements” (Bartsch and
 1601 Wellman 1995, 60). In other words, there is no single pattern that a child is
 1602 trying to master in being a competent user of “knows” and its cognates.

1603 What is interesting about this study as Baz rightly observed is that it is one
 1604 of the few scientific studies that have focussed on philosophically interesting
 1605 terms like “knows” and its cognates. Most scientific studies about words and
 1606 concepts are usually too broad in their scope and coverage to tell us what we
 1607 need to know in doing philosophy. This is important because although the
 1608 empirical result is not yet conclusive, it indicates that ordinary words like
 1609 “table” are not just like philosophically interesting words like “knows”; the

1610 latter is more complex and traces no single or simple pattern *pace* Jackson.⁵
1611 It also indicates, as Baz argued, that human language is open-ended, that is,
1612 capable of being used to make completely new moves not just at the level of
1613 the whole of an utterance but at the level of the individual parts or words in
1614 a way that is problematic for the atomistic-compositional view of language
1615 and linguistic competence. For present purposes, we can take the current
1616 empirical evidence for granted, and inquire into how to make sense of it.

1617 Baz thinks that the best way to make sense of the data is a view that combines
1618 contextualism and social-pragmatism, a view whose central ingredients
1619 come from Wittgenstein's (1953) *Philosophical Investigations* and Merleau-
1620 Ponty's (2002). Following Wittgenstein, Baz argues that we need to think of
1621 meaning as use in the sense that the significance of words depends not on
1622 their referring to items but "on whether and how we use the words, on our
1623 *meaning* them in one way or another, in a context suitable for meaning them
1624 in *that way*" (Baz 2017, 130). The advantage of the usage view in Baz's opinion
1625 is that it shows clearly that our words need not be representational and need
1626 not be thought of as naming items in the non-linguistic world to be suitable
1627 for different uses.⁶

1628 Following Merleau-Ponty, Baz argues that we need to reclaim the place
1629 of the *actual speaker* in the speech act, "the person who finding herself in
1630 some particular situation or other, may find herself moved, motivated, to
1631 speak (or think)" (Baz 2017, 131). This means that understanding the speech
1632 of another is not merely the putting together of the already fixed meaning
1633 of her words, but "coming to see her point," meaning coming to see her
1634 cares, her commitments, her history, how she sees the situation, and so on.
1635 In a significant sense therefore, the view reverses the direction of linguistic
1636 meaning implied in the atomistic compositional view: we understand the
1637 parts of speech by first understanding the whole of it, and that requires
1638 understanding the point of the actual speaker. In this connection, Baz notes
1639 that:

5 In fact, we do not need the study of how children acquire "knows" and its cognates to realise that words like "knows" do not trace a simple pattern that can be framed in terms of necessary and sufficient conditions for all instances of knowledge. We already have reasons to suspect that this is so from the failure to produce a simple account of necessary and sufficient conditions for knowledge in analytic epistemology (Shope 1983).

6 Baz says too that Wittgenstein's comparison of words to game pieces also lends credence to this idea of language.

The notion of “motive” is very important to Merleau-Ponty’s avoidance of both mechanistic and intellectualist approaches to the understanding of behavior in general and linguistic expression in particular (see (Merleau-Ponty 2002, 48–50)). On Merleau-Ponty’s way of looking at things, our speech (and behavior more generally) is normally *motivated*, in the sense that we are not merely *caused* mechanically to speak, and in the sense that our behavior manifests an *understanding* of the phenomenal world to which we respond. (Baz 2017, 131, n.14)

Baz argues that this view of language and linguistic competence gives support to a social-pragmatic account of conceptual competence inspired by Michael Tomasello (2003, 2008). On this view, in being a competent employer of “knows” and its cognates, what the child learns is different actual constructions of speech and their communicative functions, or more plainly, “stored exemplars of utterances” (Baz 2017, 162) “and what commitments (liabilities, risks) one takes upon oneself when using the words in one way or another, and in responding in one way or another to other people’s use of them” (Baz 2017, 169).

Furthermore, Baz thinks that if we accept this way of thinking about language, linguistic competence, and conceptual competence, the PMOC would be found to be seriously defective. How so? Well, if understanding the speech of another is coming to see the point of an *actual* speaker, which means coming to see her *cares*, her *commitments* vis-à-vis the question, and what *risks* and *liabilities* she may assume in answering the question one way or the other, and what empirical options we might explore to investigate whether things are thus and so, and what practical interest makes that question intelligible either to us or to the speaker, and how what is said in that context may influence what we do after; it seems clear that these conditions are lacking in the context of the PMOC. And it is because Baz thinks these conditions—let us call them “social-pragmatic conditions”—are not so realised in the PMOC that he takes the PMOC to be deeply defective and the questions asked in that context to be pointless as well. Put more generally, the view is the following:

THE SOCIAL PRAGMATIC VIEW OF LANGUAGE AND LINGUISTIC COMPETENCE. If Hearer *H* in a context *C* understands the speech of a speaker *S*, *H* does so only if the social pragmatic conditions are realised in context *C*.

1676 Notice that the view is silent as to the further question of whether the social
1677 pragmatic conditions are the only conditions required for linguistic under-
1678 standing to be possible or for words to be meaningfully used. It merely says
1679 that the social-pragmatic conditions are essential or necessary for words to
1680 do their work and for questions to have intelligibility.

1681 One urgent question is, why commit Baz to the broader goal of demarcating
1682 the region of the meaningful use of words rather than the more modest view
1683 that the questions asked by the practitioners of the PMOC are problematic
1684 or pointless?⁷ Or put differently, why think that Baz's criticism concerns the
1685 descriptions of the PMOC rather than the questions themselves and whether
1686 or not the questions are pointless? Well, the short answer is that the questions
1687 themselves are pointless precisely because the social pragmatic conditions
1688 for the felicitous use of words by both hearers and speakers are lacking in
1689 the description of the case. Baz says this precisely when he tries to show
1690 how his project fits within a broader demarcation argument that goes back to
1691 champions of experimental philosophy such as Jonathan Weinberg (2007),
1692 and more recently Edouard Machery (2017). This kind of argument relies on
1693 showing that there is a discontinuity between the scenarios described in the
1694 PMOC and the scenarios that we regularly encounter in everyday situations
1695 in a way that makes the former bad and the latter good. However, doing that
1696 often requires coming up with a set of properties defining one context but not
1697 the other context.⁸ Here is textual evidence that lends support to construing
1698 Baz in this way.

1699 The argument of this book is meant to show that the discontinuity
1700 is primarily a matter, not of *the sorts of cases* theorists have tended
1701 to focus on, as Weinberg has suggested, but of *the peculiar context*
1702 in which we attend to those cases and try to answer the theorist's
1703 questions. (Baz 2017, 33, n.33)

1704 And again:

1705 [But] if as I will argue, the ordinary and normal *conditions* for
1706 the felicitous use of the word (or concept) under investigation are
1707 lacking in the theoretical context—and, again, lacking by design—

7 An anonymous reviewer for this journal pressed me on this objection.

8 For replies to Weinberg's claim of discontinuity, see Cappelen (2012) and Nado (2015); and for a reply to Machery, see Nado (2022).

1708 then there is good reason to worry that the theorizing is bound to
1709 distort what it aims to clarify. (Baz 2017, 3)

1710 Notice that the theoretical context is also the peculiar context. Notice too that
1711 if we seek to restrict Baz's demarcation only to occasions of speech when
1712 terms like "knows" and "cause" are featured, this would be *ad hoc*. The reason
1713 these terms retain philosophical interest is due to their everyday provenance.
1714 Indeed, 'knows' and its cognates are some of the most ubiquitous terms in
1715 human language.

1716 There are two worries I would like to point out here. The first is this. Baz's
1717 claim of discontinuity implies that in the peculiar context of the PMOC, some
1718 essential conditions for the felicitous use of human words are lacking in a
1719 way that problematises the kind of questions philosophers are prone to ask
1720 in that context, as well as the answers they give. But this stands in tension
1721 with the open-endedness of human language. How so? The idea that language
1722 is open-ended, if it means anything really, means that whatever set of
1723 conditions we can identify and establish as part of the normal and felicitous
1724 use of language and words, there would always be occasions where those
1725 conditions are unmet, and yet a speaker with some ingenuity employs it in a
1726 meaningful way; a way that transmits knowledge or understanding or that
1727 serves other useful functions. Of course, language is not a human practice
1728 where anything goes. However, the thought is that given proper context, speakers
1729 and hearers can always tell the difference between what is meaningful
1730 and what is not without any predetermined criteria. Further, the thought is
1731 that these criteria, if any, would not be something that can be captured in
1732 any principled way and articulatable as something like some social pragmatic
1733 conditions. Moreover, the realm of meaning and meaningful questions and
1734 answers involving terms like "knows" and "cause" is not correctly restricted
1735 to the realm of the pragmatic or the practical for creatures like us. And that
1736 is because human beings have a capacity to engage meaningfully in things
1737 that transcend their self-interest. It seems that for evolutionary reasons, this
1738 would be a good thing. Information that has no pragmatic import for a hearer
1739 in a given context and at a particular time can have life-saving significance
1740 for that agent in a different context at a future time or perhaps for close kin.
1741 Edward Craig has a similar story of how our practically oriented concept of
1742 knowledge evolved into a more objectivised and demanding standard, where
1743 a high degree of reliability even in an improbable world is built into it. Thus,
1744 he says:

1745 In saying that someone knows whether p we are certifying him
1746 as an informant on that question, and we have no idea of the
1747 practical needs of the many people who may want to take him
1748 up on it; hence a practice develops of setting the standard very
1749 high, so that whatever turns, for them, on getting the truth about,
1750 we need not fear reproach if they follow our recommendation.
1751 (Craig 1990, 94)

1752 Perhaps it is also why “knows” and its cognates have some exceptional quali-
1753 ties such as being lexical universals, with the rare quality of being in the core
1754 vocabulary of all known human languages (Haspelmath and Tadmor 2009),
1755 and having a one-word equivalent in all natural languages (Goddard 2010).

1756 The second worry: Baz is assuming that in the theoretical or peculiar con-
1757 text of the PMOC, nothing hangs for the hearers and speakers, or the thought
1758 experimenter and his or her audience except for a theoretical interest, namely,
1759 the affirmation or the refutation of a view. But can we take that assump-
1760 tion for granted? I think not. For very often, the success of counterexamples
1761 or more generally, philosophical cases is decisive for the dominance of a
1762 particular theory and field of research. Think about the debate between com-
1763 patibilism and incompatibilism, internalism and externalism, physicalism
1764 and anti-physicalism and the decisive role that thought experiments played
1765 in those debates like Mary the colour scientist case (Jackson 1982), Gettier
1766 cases (Gettier 1963) and Truetemp case (Lehrer 1990). True enough, we only
1767 care about the truth or facts that obtain or do not obtain in those cases rather
1768 than their instrumental value. And yet because of the role those cases play
1769 in the rise and fall of certain fields of research and research prospects, it is
1770 fair to say that the facts that obtain or fail to obtain in those cases make those
1771 cases stand in the same relation to real or actual situations that are of interest
1772 to Baz: They are not idle issues to which we feel unconcerned and to which
1773 our interests, cares, and commitments are unrelated.⁹

1774 In the next section, I discuss a further challenge for Baz's account, namely,
1775 the problem of malapropism, which shows that sometimes the conditions for
1776 the ordinary use of our words are violated, and yet linguistic understanding

9 In the same vein, it is not clear that there is nothing we can do to find out whether the verdict in the cases is correct or incorrect. Indeed, this is what experimental philosophers have been doing. Although one might argue that consensus or corroboration is not correctness of intuitional judgements. But so too are perceptual judgements.

1777 is possible. This then sets the stage for presenting and developing a Davidson-
1778 inspired alternative view of language and linguistic competence.

1779 **3 The Skill or Virtue-Based Account of Language and** 1780 **Linguistic Competence**

1781 In his later writings, Davidson found the problem of “malapropism” very
1782 perplexing. Dealing with this problem led him to a view of language that
1783 affirms a continuity between linguistic competence and intellectual abilities
1784 more generally. To be sure, malapropism is a ubiquitous phenomenon in
1785 human language and registers

1786 our ability to perceive a well-formed sentence when the actual
1787 utterance was incomplete or grammatically garbled, our ability
1788 to interpret words we have never heard before, to correct slips of
1789 the tongue, or to cope with new idiolects. (Davidson 1986, 95)

1790 On the standard view of language and linguistic competence, what a hearer
1791 needs to be able to interpret a speaker is something like a complex theory or
1792 rule plus the ability to use this rule or theory or generalisation in a systematic
1793 way to make sense of novel situations. Further, because this capacity is taken
1794 as a learned convention, one that is shared between hearers and speakers,
1795 it is something that the hearer has in advance of the occasion of linguistic
1796 exchange. Notice that this standard view is also the view defended by Jackson
1797 as previously presented and discussed (Jackson 2011). Recall that on that view,
1798 namely, the atomistic-compositional view, language is like the numbering
1799 system where there are finite numerals that can be used to generate complex
1800 ones infinitely. Speakers and hearers have this system in advance of particular
1801 linguistic exchanges.

1802 However, the phenomenon of malapropism challenges this notion because
1803 the competence (or capacity) that it calls for from the hearer is not part of
1804 what normally constitutes one’s basic linguistic competence, mastered in
1805 advance of the occasion of linguistic exchange. Indeed, as Davidson points
1806 out, the fact that makes the theory or rule general equally makes it unsuitable
1807 to cope with the particular linguistic habits of different individuals, say that
1808 of Mrs. Malaprop’s “nice derangement of epitaphs” being “nice arrangement
1809 of epithets”.¹⁰ More generally, the theory or rule is unhelpful in coping with a

10 Malaprop was a character famous for her verbal blunders in Richard Sheridan’s play *The Rivals*.

1810 particular speaker at a particular time in a particular occasion. This applies to
1811 Baz's account too since for him there are "ordinary and normal conditions for
1812 the felicitous use" of human words or concepts (Baz 2017, 3), conditions which
1813 he thinks are lacking in the context of the PMOC. But then, in malapropism
1814 such as grammatically garbled utterances and slips of tongues, those normal
1815 conditions for the felicitous use of words and for their "functioning as they
1816 do" in ordinary discourse (Baz 2017, 22) are violated. Further, it is not the case
1817 that for Baz there is one generic condition, namely, that one's utterance has a
1818 point. On the contrary, that one's utterance has a point is fixed by it satisfying
1819 "the ordinary and normal conditions" for the felicitous use of human words
1820 and for meaning words one way or the other. For he says:

1821 And the basic problem with so much philosophizing, both tradi-
1822 tional and contemporary—the basic problem with the method of
1823 cases as commonly practiced, for example—is that the philoso-
1824 pher either takes his words to mean something clear even apart
1825 from *his* meaning something clear by means of them, or else takes
1826 himself to be able to mean his words in some determinate way,
1827 *even though the conditions for thus meaning his words are missing*
1828 *in his particular context* and cannot be created by a sheer act of
1829 will, or by concentrating one's mind in some special way. (Baz
1830 2017, 141, italics mine)

1831 Here is an additional challenge from malapropism to any generic view
1832 of language and linguistic competence. Sometimes in linguistic exchange,
1833 linguistic understanding is transmitted despite the hearer completely mistak-
1834 ing the speaker's verbal communication and vice versa. Davidson gives an
1835 example of such a case:

1836 When I first read Singer's piece on Goodman Ace, I thought that
1837 the word 'malaprop', though the name of Sheridan's character, was
1838 not a common noun that could be used in place of 'malapropism'.
1839 It turned out to be my mistake. Not that it mattered: I knew what
1840 Singer meant, even though I was in error about the word; I would
1841 have taken his meaning in the same way if he had been in error
1842 instead of me. We could both have been wrong, and things would
1843 have gone as smoothly. (Davidson 1986, 90)

1844 Here as elsewhere, learned convention breaks down and the conditions for
1845 the normal and felicitous use of words are violated and yet linguistic under-

1846 standing is transmitted or made possible. The question is, how is this possible?
 1847 What capacity does the hearer (and speaker) depend on? Davidson makes the
 1848 following suggestion:

1849 This characterisation of linguistic ability is so nearly circular that
 1850 it cannot be wrong: it comes to saying that the ability to communi-
 1851 cate by speech consists in the ability to make oneself understood,
 1852 and to understand. It is only when we look at the structure of this
 1853 ability that we realise how far we have drifted from standard ideas
 1854 of language mastery. For we have discovered no learnable com-
 1855 mon core of consistent behaviour, no shared grammar or rules,
 1856 no *portable* interpreting machine set to grind out the meaning of
 1857 an arbitrary utterance. We may say that linguistic ability is the
 1858 ability to converge on a passing theory from time to time—this is
 1859 what I have suggested, and I have no better proposal. But if we
 1860 do say this, then we should realise that we have abandoned not
 1861 only the ordinary notion of a language, but we have erased the
 1862 boundary between knowing a language and knowing our way
 1863 around in the world generally. (Davidson 1986, 445–446, italics
 1864 mine)

1865 We can summarise the import of this account as follows:

1866 THE SKILL OR VIRTUE-BASED ACCOUNT OF LANGUAGE AND LIN-
 1867 GUISTIC COMPETENCE. If Hearer *H* in a context *C* understands the
 1868 speech of a speaker *S*, *H* does so in virtue of her skills or virtues.

1869 The rationale for speaking of skills or virtues here is two-fold. First, it is to
 1870 pick up on a suggestion by Davidson when he talks about the skillful hearer
 1871 (and speaker) as being one that can get along well in linguistic exchanges and
 1872 performances without needing mastery or knowledge of Gricean principles,
 1873 because these general principles “are a kind of *skill* we expect of an interpreter
 1874 and without which communication would be greatly impoverished” (David-
 1875 son 1986, 437). Relatedly, he talks about virtues such as practical wisdom,
 1876 intelligence, and wit as the non-linguistic competencies we rely on in getting
 1877 things right from time to time, occasion to occasion (Davidson 1986, 446).
 1878 Davidson also mentions luck. But here luck is not a capacity of speakers or
 1879 hearers. Rather, it merely refers to their being in a favourable environment
 1880 such that under normal circumstances, when they attempt to understand one

1881 another in linguistic exchange, they achieve that aim. Further, I persist in
1882 speaking of “skills and virtues” because although all skills can be classified
1883 as virtues of agents, not all virtues can be classified as skills. One particular
1884 exception to this is practical wisdom (Stichter 2018). Let us take these points
1885 in turn. First, virtues are skills because acting well is much like working well
1886 (Annas 1995) and both involve practices of self-regulation to achieve a goal:
1887 in one case, the goal of acting well, and in the other case, the goal of working
1888 well (Stichter 2018). And second, although practical wisdom involves some
1889 elements of skills, namely, making good judgements in particular situations,
1890 it also involves other dimensions, namely, considering how one’s action fits
1891 into an overall conception of the good life (Stichter 2018). So, while it might
1892 be true that agents rely on some aspects of practical wisdom in order to act
1893 well in particular situations and to get along in a linguistic exchange, practical
1894 wisdom in itself is too broad and varied to be classified merely as a set of skills.

1895 Furthermore, the competent hearer (and speaker) would also recruit other
1896 capacities of the virtuous agents. Of particular importance in the present
1897 context would be “sensibility.” In her discussion of the virtues (and the vices
1898 of the mind), Alessandra Tanesini defines sensibility as a disposition to “use
1899 one’s perceptual capacities in distinctive ways in the service of epistemic
1900 activities” (Tanesini 2021, 27). The example she gives is the observant person:

1901 The person who is observant has reliable vision but he also experi-
1902 ences as salient those features of the visual field that are relevant
1903 to his epistemic aims. He directs visual attention to these aspects of
1904 the environment. By directing attention to them, and thus putting
1905 them at the centre of his visual field, he is able to take in more
1906 detail about these items since foveal vision has a higher degree
1907 of resolution than peripheral vision. Had those items remained
1908 at the periphery of his vision, many of their features would have
1909 remained undetected. If this is right, being observant is the com-
1910 plex disposition to detect the salient aspects of the environment
1911 by experiencing feelings that direct one’s attention towards these
1912 features. (Tanesini 2021, 27–28)

1913 Applied as a competence essential to linguistic understanding, sensibility is
1914 an auxiliary competence, an enabler of visual and auditory competencies of
1915 agents. And what that means precisely is that it makes it possible for one to put
1916 to use those primary competencies in picking up what is being passed across,

1917 verbally and non-verbally, where this is something that can be missed easily
 1918 if one is not attentive to another's peculiar linguistic habits in the context of
 1919 linguistic exchange.

1920 The second rationale for the skill, or virtue-based model, is that it allows
 1921 us to cash out the Davidson-inspired view in a way that makes the relevant
 1922 competence an instance of a more general and familiar kind of know-how.
 1923 One difficulty that we can resolve in Davidson's account if we take seriously
 1924 the virtue or skill-based model is how to understand a practice that is non-rule-
 1925 based and yet rational and well-ordered. And the thing to say is that in both
 1926 virtue and skills, we already have human practices that are well-regulated
 1927 without the agents relying on rules. Take the skill-based model. Following
 1928 this model, I am suggesting that knowing a language is much like knowing
 1929 how to drive a car. In the beginning, the driver learns rules of thumb such as
 1930 "shift up when the motor sounds like it is racing and down when it sounds
 1931 like it is straining."¹¹ As Dreyfus and Dreyfus who have studied human skills
 1932 in various domains of performance argued:

1933 It seems that beginners make judgements using strict rules and
 1934 features, but that with talent and a great deal of involved experi-
 1935 ence the beginner develops into an expert who sees intuitively
 1936 what to do without applying rules and making judgements at all.
 1937 (Dreyfus and Dreyfus 1991, 235)

1938 On this thinking, if one is following rules in a practice, that just shows one
 1939 is not yet proficient in that practice. The same story applies to the virtuous
 1940 agent. As Linda Zagzebski puts it: "Persons with practical wisdom learn how
 1941 and when to trust certain feelings, and they develop habits of attitude and
 1942 feeling that enable them to reliably make good judgments without being
 1943 aware of following a procedure" (Zagzebski 1996, 226). Notice too the role of
 1944 the virtues and skills here: they are dispositions that allow agents to act in a
 1945 systematic and organised way and to do so well in a context where the relevant
 1946 practice is not rule-governed. Plausibly, the reason this is so is because both
 1947 skill and virtues have a kind of *logos*, in the sense that they have an intrinsic
 1948 intellectual structure built into them (Bloomfield 2000). Mastering a skill,
 1949 including language, is mastering this *logos*; and thus, possessing the practical

11 Such rules of thumb are just heuristics or generalisations about language that hold for the most part.

1950 intelligence to act and to sensibly follow the actions of others and to solve
1951 problems in the relevant domain or activity.

1952 From this standpoint, we can appreciate another respect in which the skill
1953 or virtue-based account and Baz's view diverge. On Baz's account, the motive
1954 of the speaker plays an essential role in coming to see the point of the speaker.
1955 Notice that "motive" here does not mean intention. It means rather the "moti-
1956 vating factors", which are internal to the perspective of the speaker, namely,
1957 the cares, the commitments, the risks and the liabilities of the speaker. On the
1958 other hand, for the skill or virtue-based account, that component is not always
1959 essential even though it sometimes can form a part of the process of coming
1960 to see the point of the speaker's utterance. Indeed, I believe that that form
1961 of internalism about linguistic sense, or meaning, was part of the tradition
1962 of thought that Gilbert Ryle tries to wean analytic philosophy from (see also,
1963 [Putnam 1975b](#)) when he argued that we should think of understanding as
1964 *knowing how* and linguistic understanding including, as an exercise of that
1965 *knowing how*. He writes:

1966 Understanding a person's deeds and words is not, therefore, any
1967 kind of problematic divination of occult processes. For this divi-
1968 nation does not and cannot occur, whereas understanding does
1969 occur. Of course, it is part of my general thesis that the supposed
1970 occult processes are themselves mythical; there exists nothing
1971 to be the object of the postulated diagnoses. But for the present
1972 purpose it is enough to prove that, if there were such inner states
1973 and operations, one person would not be able to make probable
1974 inferences to their occurrence in the inner life of another. ([Ryle](#)
1975 [2009, 41](#))

1976 Let me elaborate more on what this rejection of the internalistic picture in
1977 the motivating sense means by commenting on what Ryle is getting at here.
1978 Suppose I am playing chess with Magnus Carlsen, the Norwegian grandmaster.
1979 He makes a particular opening move that seems initially surprising to me.
1980 But as a fellow grandmaster who is equally skilful or competent in the game
1981 and who has sufficient experience dealing with a move like that, I can know
1982 what that move is about without caring about what has made Carlsen make
1983 this move. I can know that a move like that in a context like this means that a
1984 particular form of attack on my king is imminent and that moving my pieces
1985 in a specified way is the best way to counter it. The same is true of "moves" in

1986 linguistic performances, as Baz would like to call human utterances or the use
 1987 of words in language. Hearers can tell that an utterance like this in a context
 1988 like that means so and so without caring about what has moved the speaker
 1989 to say so and so.

1990 With this view of language and linguistic competence in mind, let us address
 1991 two challenges in connection with the PMOC. The first challenge here is to
 1992 explain how, as competent speakers, we are able to understand and answer
 1993 the questions that philosophers often ask in the context of the PMOC, such
 1994 as, does the protagonist in that scenario know so and so? And the second
 1995 challenge is how to make the aim of using the PMOC intelligible in the light
 1996 of the complexity of human language, that is, without glossing over that very
 1997 complexity. I take each in turn.

1998 On the skill or virtue-based view, competent speakers can understand and
 1999 answer the questions of the sort “does *X* know *Y*?” not because they have
 2000 latched onto the pattern of “knows” *pace* Jackson or because they possess
 2001 stored exemplars of utterances and knowledge of the communicative motives
 2002 of speakers *pace* Baz. On the contrary—when they do, that is in virtue of their
 2003 having mastered a technique in the use of “knows” and its cognates. In fact,
 2004 this suggestion finds its earliest expression in the later Wittgenstein when he
 2005 says:

2006 The grammar of the word “know” is evidently closely related
 2007 to the grammar of the words “can”, “is able to.” But also closely
 2008 related to that of the word “understand” (*To have ‘mastered’ a*
 2009 *technique*). (Wittgenstein 1953, sec. 150–151, italics mine)¹²

-
- 12 Should we read Wittgenstein’s suggestion as the mastery of grammatical rules or relationships? There is abundant evidence in the text and elsewhere that that is not what Wittgenstein had in mind. To start with, in the paragraphs that followed this statement (i.e., Wittgenstein 1953, sec. 151–152), he says that it is conceivable that the relevant formula (or rule or grammatical relationship) might occur to the speaker and yet the speaker fails to understand. Further, in an unpublished manuscript, translated by Norman Malcolm (1989), Wittgenstein writes: “Often one can say: this pattern looked at so, must have this continuation. I want, however, to stipulate an ‘interpretation’ [*Auffassung*], (something like the old ‘Proposition’), which determines the series like an infallible machine through which a conveyor belt runs. So that only this continuation fits this interpretation. In reality, however, there are not two things that here fit together. But one can say: You are by your training, so adjusted [*eingestellt*], that always, without reflection, you declare some definite thing to be that which fits. Something that agrees with what others declare to be what fits” (Wittgenstein, Unpublished manuscript, 86–87; cited in Malcolm (1989)). On this view, it is by one’s training as a member of a shared community and practice that one is able to reliably employ terms like “knows” and extend the practice in similar situations. For Wittgenstein, that

2010 Such skills or techniques are suitably grounded in experience in such a way
2011 that the agents exercising them can always be counted upon to answer such
2012 questions in a range of situations, not only in actual ones but in possible ones
2013 that bear similarity to the actual ones, where what is “similar” cannot be
2014 established in any rigid way, for example, through the claim of discontinuity
2015 between the context of the PMOC and everyday contexts. Indeed, as argued
2016 earlier, being competent users of “knows” and answering questions such as
2017 “does X know Y?” in a range of situations might be part of our evolutionary
2018 heritage. Also, a recent trend in cognitive science seems to lend support
2019 to this skill-based suggestion. Here is Lawrence Barsalou and colleagues
2020 summarising the emerging consensus here:

2021 [C]onceptual knowledge is not a global description of a category
2022 that functions as a detached database about its instances. Instead,
2023 conceptual knowledge is the *ability* to construct situated con-
2024 ceptualizations of the category that serves agents in particular
2025 situations. (Barsalou et al. 2003, 89)¹³

picture of a skill or technique grounded in training replaces the picture of the grammatical rule acting like an infallible conveyor belt that determines its extension in novel situations.

- 13 As previously pointed out, Baz argues that the atomistic-compositional view seems to go side by side with the assumption that the primary purpose of language is to transmit information, that is, it seems to go side by side with the representational-referential view of language. Again, there is no need to hold on to that problematic assumption. All that is necessary for the philosophical method of cases to get off the ground once the atomistic-compositional view is set aside and the skill or virtue-based view is assumed is that among other things, language can be used to transmit information, where again given appropriate context agents can tell when this is the case. In fact, the empirical study that Baz analyses in support of his view does not presuppose otherwise. To see this, notice that although in Baz's discussion of this study, he cites the frequency with which children refer to their own mental states as clear vindication of his view of language, the data also show that this frequency diminishes as the children grow older. Bartsch and Wellman also note that “our data provide no evidence that a representational understanding of beliefs is a significantly later achievement, following only on the heels of an earlier ‘connections’ misconstrual of beliefs” (Bartsch and Wellman 1995, 57). Further, even in their first-person reference to mental states, the data do not contradict representational presuppositions. As the authors put it “[W]hen children first use *know* to refer to people's knowledge in our data, in their utterances coded as genuine psychological references, they primarily refer either to situations involving successful actions or to correct statements” (Bartsch and Wellman 1995, 60). And lastly, in an earlier study of our everyday conception of knowledge as manifested in words like “knows” and “knew”, Perner (1991) shows that knowledge is associated with success and successful actions, with factual states of affairs and is formed by exposure to the relevant information or experience.

2026 Now the second challenge. Using the PMOC, Edmund Gettier drew the
2027 attention of the philosophical community to an aspect of knowledge, namely,
2028 that the term is a success notion; the term does not apply to someone whose
2029 belief is chancy or accidental. Does that gloss over the complexities in our use
2030 of “knows” and its cognates? Baz thinks so (see Baz 2017, 122). But there are
2031 good reasons to doubt that conclusion. To start with, notice that the idea that
2032 knowledge is a success term is implied in the result of the study of Bartsch
2033 and Wellman (1995). Further, imagine as we do in the analysis of knowing
2034 that we highlight “success” or “achievement” as a salient feature of the term
2035 “knows” and explain knowledge in terms of these notions (Greco 2010). I
2036 argue that doing so does not obscure the subject matter of philosophy as
2037 Baz implies. On the contrary, doing so advances our understanding of the
2038 subject matter. Indeed, this is closely related to scientific practice. Biologists
2039 know that the term “fish” picks out various kinds of properties such as having
2040 fins, having scales, having a tail, breathing underwater, being oviparous, not
2041 suckling one’s young, and being cold-blooded. But from the point of view of
2042 understanding, and classifying future unknown cases, they merely highlight
2043 a fewer set of properties rather than all of the above, especially those that are
2044 natural and explanatory so that the term “fish” is used to refer to a completely
2045 aquatic, water-breathing, cold-blooded craniate vertebrate (Slote 1966). I be-
2046 lieve the same story applies here to the PMOC in the analysis of knowing. In
2047 highlighting the fact that knowledge is a success term, we are able to track
2048 something important, deep and explanatory about this phenomenon, some-
2049 thing we can also use to understand other terms or concepts or issues. For
2050 example, knowledge firsters use the suggestion that knowledge is a success
2051 term to understand the notion of intellectual ability or competence (Kelp
2052 2021).

2053 Let us conclude this section by noting how the skill or virtue-based model
2054 of language and linguistic competence shares something positive with Baz’s
2055 social pragmatic account. Clearly, both recover the place of the speaking
2056 subject and reject the idea implied in the atomistic-compositional view that
2057 human words can speak for themselves, “over our heads as it were—and of
2058 language as a system of significant signs that does not depend on speakers (and
2059 listeners) for its ongoing maintenance” (Baz 2017, 96). Indeed, in evaluating
2060 Gettier cases, for example, we often need to tell whether or not and in what
2061 relevant sense the cases we are evaluating resemble clear instances where
2062 the property or term is clearly instantiated in a case. And “which way one
2063 goes depends on what one finds normal or natural, which partly depends

2064 on the past course of one's sense experience" (Williamson 2007, 190). Notice
2065 that the capacity to tell that something is "normal or natural" is much in
2066 line with the capacity that comes with practical wisdom, which is shaped by
2067 experience, including sense experience, and expressed in habits of attitude
2068 and feeling that enable one to reliably make good judgements without being
2069 aware of following any rule. Moreover, in a non-actual instance of a Gettier
2070 case, readers often need to follow "in their own imaginative construction the
2071 lead of the author of the examples" (Sosa 2009, 107), and they have to fill
2072 out the details of the stories, which are often partial and incomplete. Here as
2073 elsewhere too, one needs to tell whether or not and in what relevant sense the
2074 case one is evaluating resembles clear instances where the property or term
2075 is clearly instantiated. Moreover, which way one goes depends on what one
2076 finds normal or natural. Notice also that if the kind of story that particularists
2077 such as Jonathan Dancy tell about the use of thought experiments in moral
2078 philosophy is true, namely, that no suitable supply of general principles can
2079 help the moral agent in picking out what is morally salient about a case
2080 (Dancy 1985), then we have good reason to believe that even here what the
2081 agent does is to recruit the kind of capacities that the skill or virtue-based
2082 model highlights. In any case, a theory of language and linguistic competence
2083 begins from the correct assumption that ordinary speakers already do well in
2084 linguistic performances and presents an explanation of how speakers are able
2085 to so perform. I have argued that once we reject the atomistic compositional
2086 view, it does not follow that we must embrace the social pragmatic story and
2087 all the problems it poses for the PMOC.¹⁴

14 An anonymous reviewer for this journal pressed the following worries. The first worry is that "the proposed virtue-based account of linguistic understanding is perfectly compatible with there being cases/situations in which it doesn't make sense to ask about a certain subject and a certain fact 'Does *S* know that *p*?' Hence, it seems to me that further argument is needed in order to make the case for the meaningfulness of the theorist's questions about the philosophical thought experiments discussed in Baz." Reply: The worry that there are particular cases, say some very outlandish cases, where it does not make sense to ask about a certain subject and a certain fact 'Does *S* know that *p*?' does not licence the general or global worry about the PMOC as discussed in Baz. Even mainstream philosophers themselves have expressed concern that some cases are so outlandish that they are not theoretically useful because they do not resemble cases we face in everyday life (Weatherson 2003, 8). Here is another related worry pressed by the reviewer: "Davidson will also need some distinction (or demarcation) between situations in which the utterances of a certain sentence, e.g., of the form '*x* knows that *p*' makes sense and situations in which it doesn't (because obviously, you cannot meaningfully utter just any sentence in any context). And it is not obvious to me that according to Davidson the first kind of situations won't be exactly the ones in which the relevant utterance has a point." Reply: It is not exactly clear why

2084 4 Conclusion

2089 In this paper, I have argued essentially that the philosophical method of cases
 2090 does not need to presuppose the problematic view of language and linguistic
 2091 competence Baz attributes to its practitioners or defenders—the atomistic
 2092 compositional view. And neither do friends of the PMOC need to embrace the
 2093 social pragmatic view that Baz presents with all its negative consequences for
 2094 the PMOC. Let me end with where the Davidson-inspired skill or virtue-based
 2095 view leaves us in terms of the epistemology of philosophy. In my opinion, it
 2096 lends independent support to the view, now current in the epistemology of
 2097 philosophy that the epistemology of philosophy is an application of social
 2098 epistemology. Williamson (2007); Nagel (2012) and more recently Irikefe
 2099 (2022) champion this epistemological thesis and it seems to me the right
 2100 way to explain how philosophical knowledge is possible and how it can be
 2101 defended against various challenges posed against it.*

2102 Paul O. Irikefe

2103  0000-0003-4086-8229

2104 University of California, Irvine

2105 Irikefef@uci.edu

2106 References

- 2107 ANNAS, Julia. 1995. "Virtue as a Skill." *International Journal of Philosophical Studies*
 2108 3(2): 227–243, doi:10.1080/09672559508570812.
- 2109 BARSALOU, Lawrence W., SIMMONS, W. Kyle, BARBEY, Aron K. and WILSON, Christine
 2110 D. 2003. "Grounding Conceptual Knowledge in Modality-Specific Systems." *Trends*
 2111 *in Cognitive Sciences* 7(2): 84–91, doi:10.1016/S1364-6613(02)00029-3.

we need such a demarcation to start with since we can get along without it, and if a theory of language is an attempt to model what people already do, that is exactly the kind of story our theory should also be telling. Further, one might wonder whether such a demarcation does not imply by its very existence that there is a rigid boundary of what counts as meaningful linguistic occasions and what does not count as so. It does in my opinion. And it is why in the history of philosophy these kinds of projects, which seek to demarcate some regions of language as linguistically acceptable and others that are not on the basis of some criteria have had little or no success. In any case, the Davidson-inspired view shows us a way to proceed without it.

* I am grateful to two anonymous reviewers for this journal for their invaluable comments that helped shape the development of the paper to its current form. In the same vein, I thank Alessandra Tanesini and Mona Simion. Their comments were helpful in refining the paper's central claims.

- 2112 BARTSCH, Karen and WELLMAN, Henry M. 1995. *Children Talk About the Mind*. Oxford:
2113 Oxford University Press, doi:10.1093/oso/9780195080056.001.0001.
- 2114 BAZ, Avner. 2016. "Recent Attempts to Defend the Philosophical Method of Cases
2115 and the Linguistic (Re)turn." *Philosophy and Phenomenological Research* 92(1):
2116 105–130, doi:10.1111/phpr.12106.
- 2117 —. 2017. *The Crisis of Method in Contemporary Analytic Philosophy*. Oxford: Oxford
2118 University Press, doi:10.1093/oso/9780198801887.001.0001.
- 2119 BLOOMFIELD, Paul. 2000. "Virtue Epistemology and the Epistemology of Virtue."
2120 *Philosophy and Phenomenological Research* 60(1): 23–43, doi:10.2307/2653426.
- 2121 BOH, Ivan. 1985. "Belief, Justification, and Knowledge: Some Late-Medieval Epistemic
2122 Concerns." *Quidditas: Online Journal of the Rocky Mountain Medieval and Renais-
2123 sance Association* 6(1): 8, <https://scholarsarchive.byu.edu/rmmra/vol6/iss1/8>.
- 2124 CAPPELEN, Herman. 2012. *Philosophy without Intuitions*. Oxford: Oxford University
2125 Press, doi:10.1093/acprof:oso/9780199644865.001.0001.
- 2126 —. 2014. "X-Phi Without Intuitions?" in *Intuitions*, edited by Anthony Robert BOOTH
2127 and Darrell P. ROWBOTTOM, pp. 269–286. Oxford: Oxford University Press, doi:10
2128 .1093/acprof:oso/9780199609192.003.0015.
- 2129 CRAIG, Edward J. 1990. *Knowledge and the State of Nature: An Essay in Conceptual
2130 Synthesis*. Oxford: Oxford University Press, doi:10.1093/0198238797.001.0001.
- 2131 DANCY, Jonathan. 1985. "The Role of Imaginary Cases in Ethics." *Pacific Philosophical
2132 Quarterly* 66(1/2): 141–153, doi:10.1111/j.1468-0114.1985.tb00246.x.
- 2133 DAVIDSON, Donald. 1986. "A Nice Derangement of Epitaphs." in *Truth and Interpreta-
2134 tion: Perspectives on the Philosophy of Donald Davidson*, edited by Ernest LEPORE,
2135 pp. 433–446. Oxford: Blackwell Publishers. Also published in Grandy and Warner
2136 (1986, 157–174).
- 2137 DEUTSCH, Max Emil. 2015. *The Myth of the Intuitive: Experimental Philosophy and
2138 Philosophical Method*. Cambridge, Massachusetts: The MIT Press, doi:10.7551/mi
2139 tpress/9780262028950.001.0001.
- 2140 DREYFUS, Hubert L. and DREYFUS, Stuart E. 1991. "Towards a Phenomenology of
2141 Ethical Expertise." *Human Studies* 14(4): 229–250, doi:10.1007/bf02205607.
- 2142 EARLENBAUGH, Joshua and MOLYNEUX, Bernard. 2009. "If Intuitions Must Be Eviden-
2143 tial then Philosophy is in Big Trouble." *Studia Philosophica Estonica* 2(2): 35–53,
2144 doi:10.12697/spe.2009.2.2.03.
- 2145 ELBOURNE, Paul. 2011. *Meaning: A Slim Guide to Semantics*. Oxford: Oxford University
2146 Press.
- 2147 GETTIER, Edmund L., III. 1963. "Is Justified True Belief Knowledge?" *Analysis* 23(6):
2148 121–123, doi:10.1093/analys/23.6.121.
- 2149 GODDARD, Cliff. 2010. "Universals and Variation in the Lexicon of Mental State Con-
2150 cepts." in *Words and the Mind: How Words Capture Human Experience*, edited by
2151 Barbara C. MALT and Phillip WOLFF, pp. 72–92. New York: Oxford University
2152 Press, doi:10.1093/acprof:oso/9780195311129.003.0005.

- 2153 GRANDY, Richard E. and WARNER, Richard, eds. 1986. *Philosophical Grounds of Rationality: Intentions, Categories, Ends*. Oxford: Oxford University Press.
- 2154
- 2155 GRECO, John. 2010. *Achieving Knowledge: A Virtue-Theoretic Account of Epistemic Normativity*. Cambridge: Cambridge University Press, doi:10.1017/cbo978051184
- 2156 4645.
- 2157
- 2158 HASPELMATH, Martin and TADMOR, Uri. 2009. "The Loanword Typology Project and the World Loanword Database." in *Loanwords in the World's Languages: A Comparative Handbook*, edited by Martin HASPELMATH and Uri TADMOR, pp.
- 2159 1–34. Berlin: de Gruyter Mouton, doi:10.1515/9783110218442.1.
- 2160
- 2161 ICHIKAWA, Jonathan Jenkins and JARVIS, Benjamin W. 2009. "Thought-Experiment Intuitions and Truth in Fiction." *Philosophical Studies* 142(2): 221–246, doi:10.100
- 2162 7/s11098-007-9184-y.
- 2163
- 2164 IRIKEFE, Paul Oghenovo. 2020. "A Fresh Look at the Expertise Reply to the Variation Problem." *Philosophical Psychology* 33(6): 840–867, doi:10.1080/09515089.2020.17
- 2165 61541.
- 2166
- 2167 —. 2022. "The Epistemology of Thought Experiments without Exceptionalist Ingredients." *Synthese* 200(3): 191, doi:10.1007/s11229-022-03690-2.
- 2168
- 2169 JACKSON, Frank. 1982. "Ephiphenomenal Qualia." *The Philosophical Quarterly* 32(127): 127–136. Reprinted in Jackson (1998, 57–69), doi:10.2307/2960077.
- 2170
- 2171 —. 1998. *Mind, Method, and Conditionals: Selected Papers*. London: Routledge, doi:10
- 2172 .4324/9780203271308.
- 2173
- 2174 —. 2011. "On Gettier Holdouts." *Mind and Language* 26(4): 468–481, doi:10.1111/j.14
- 2175 68-0017.2011.01427.x.
- 2176
- 2177 KELP, Christoph. 2021. "Theory of Inquiry." *Philosophy and Phenomenological Research* 103(2): 359–384, doi:10.1111/phpr.12719.
- 2178
- 2179 LEHRER, Keith. 1990. *Theory of Knowledge*. 1st ed. Dimensions of Philosophy Series. Boulder, Colorado: Westview Press. Second edition: Lehrer (2000), doi:10.4324/97
- 2180 80203870198.
- 2181
- 2182 —. 2000. *Theory of Knowledge*. 2nd ed. Dimensions of Philosophy Series. Boulder, Colorado: Westview Press. First edition: Lehrer (1990), doi:10.4324/978042949426
- 2183 0.
- 2184
- 2185 MACHERY, Edouard. 2017. *Philosophy Within Its Proper Bounds*. Oxford: Oxford University Press, doi:10.1093/oso/9780198807520.001.0001.
- 2186
- 2187 MALCOLM, Norman. 1989. "Wittgenstein on Language and Rules." *Philosophy* 64(247): 5–28, doi:10.1017/s0031819100044004.
- 2188
- 2189 MERLEAU-PONTY, Maurice. 1945. *Phénoménologie de la perception*. Tel. Paris: Gallimard.
- 2190
- 2191 —. 2002. *Phenomenology of Perception*. 2nd ed. London: Routledge. Translation of Merleau-Ponty (1945) by Colin Smith, doi:10.4324/9780203994610.
- 2192
- 2193 NADO, Jennifer. 2015. "Intuition, Philosophical Theorizing, and the Threat of Skepticism." in *Experimental Philosophy, Rationalism, and Naturalism: Rethinking*

- 194 *Philosophical Method*, edited by Eugen FISCHER and John COLLINS, pp. 204–221.
195 London: Routledge, doi:[10.4324/9781315714196](https://doi.org/10.4324/9781315714196).
- 196 —. 2022. “Philosophizing Out of Bounds.” *Philosophical Studies* 179(1): 319–327,
197 doi:[10.1007/s11098-020-01582-0](https://doi.org/10.1007/s11098-020-01582-0).
- 198 NAGEL, Jennifer. 2012. “Intuitions and Experiments: A Defense of the Case Method in
199 Epistemology.” *Philosophy and Phenomenological Research* 85(3): 495–527, doi:[10](https://doi.org/10.1111/j.1933-1592.2012.00634.x)
200 [.1111/j.1933-1592.2012.00634.x](https://doi.org/10.1111/j.1933-1592.2012.00634.x).
- 201 PERNER, Josef. 1991. *Understanding the Representational Mind*. Cambridge, Mas-
202 sachusetts: The MIT Press, doi:[10.7551/mitpress/6988.001.0001](https://doi.org/10.7551/mitpress/6988.001.0001).
- 203 PUTNAM, Hilary. 1975a. *Mind, Language and Reality: Philosophical Papers, Volume 2*.
204 Cambridge: Cambridge University Press, doi:[10.1017/cbo9780511625251](https://doi.org/10.1017/cbo9780511625251).
- 205 —. 1975b. “The Meaning of ‘Meaning’.” in *Minnesota Studies in the Philosophy of*
206 *Science, Volume VII: Language, Mind, and Knowledge*, edited by Keith GUNDERSON,
207 pp. 131–193. Minneapolis, Minnesota: University of Minnesota Press. Reprinted
208 in Putnam (1975a, 215–271).
- 209 RYLE, Gilbert. 2009. *The Concept of Mind*. London: Routledge, doi:[10.4324/97802038](https://doi.org/10.4324/9780203875858)
210 [75858](https://doi.org/10.4324/9780203875858).
- 211 SHOPE, Robert K. 1983. *The Analysis of Knowing: A Decade of Research*. Princeton,
212 New Jersey: Princeton University Press, doi:[10.1515/9781400886555](https://doi.org/10.1515/9781400886555).
- 213 SLOTE, Michael Anthony. 1966. “The Theory of Important Criteria.” *The Journal of*
214 *Philosophy* 63(8): 211–224, doi:[10.2307/2023977](https://doi.org/10.2307/2023977).
- 215 SOSA, Ernest. 2007. “Experimental Philosophy and Philosophical Intuition.” *Philo-*
216 *sophical Studies* 132(1): 99–107, doi:[10.1007/s11098-006-9050-3](https://doi.org/10.1007/s11098-006-9050-3).
- 217 —. 2009. “A Defense of the Use of Intuitions in Philosophy.” in *Stich and His Critics*,
218 edited by Dominic MURPHY and Michael A. BISHOP, pp. 101–112. Philosophers
219 and Their Critics n. 31. Malden, Massachusetts: Wiley-Blackwell, doi:[10.1002/97](https://doi.org/10.1002/9781444308709.ch6)
220 [81444308709.ch6](https://doi.org/10.1002/9781444308709.ch6).
- 221 STICHTER, Matt. 2018. *The Skillfulness of Virtue: Improving our Moral and Epistemic*
222 *Lives*. Cambridge: Cambridge University Press, doi:[10.1017/9781108691970](https://doi.org/10.1017/9781108691970).
- 223 TANESINI, Alessandra. 2021. *The Mismeasure of the Self: A Study in Vice Epistemology*.
224 Oxford: Oxford University Press, doi:[10.1093/oso/9780198858836.001.0001](https://doi.org/10.1093/oso/9780198858836.001.0001).
- 225 TOMASELLO, Michael. 2003. *Constructing a Language: A Usage-Based Theory of Lan-*
226 *guage Acquisition*. Cambridge, Massachusetts: Harvard University Press, doi:[10.2](https://doi.org/10.2307/j.ctv26070v8)
227 [307/j.ctv26070v8](https://doi.org/10.2307/j.ctv26070v8).
- 228 —. 2008. *Origins of Human Communication*. Cambridge, Massachusetts: The MIT
229 Press, doi:[10.7551/mitpress/7551.001.0001](https://doi.org/10.7551/mitpress/7551.001.0001).
- 230 WEATHERSON, Brian. 2003. “What Good Are Counterexamples?” *Philosophical Studies*
231 115(1): 1–31, doi:[10.1023/a:1024961917413](https://doi.org/10.1023/a:1024961917413).
- 232 WEINBERG, Jonathan M. 2007. “How to Challenge Intuitions Empirically Without
233 Risking Skepticism.” in *Midwest Studies in Philosophy 31: Philosophy and the*
234 *Empirical*, edited by Peter A. FRENCH and Howard K. WETTSTEIN, pp. 318–343.

- 2235 Boston, Massachusetts: Blackwell Publishers, doi:10.1111/j.1475-4975.2007.00157.
2236 X.
- 2237 WILLIAMSON, Timothy. 2007. *The Philosophy of Philosophy*. Oxford: Blackwell Pub-
2238 lishers. Second edition: Williamson (2022), doi:10.1002/9780470696675.
- 2239 —. 2022. *The Philosophy of Philosophy*. 2nd ed. The Blackwell / Brown Lectures in
2240 Philosophy. Hoboken, New Jersey: Wiley-Blackwell. First edition: Williamson
2241 (2007), doi:10.1002/9781119616702.
- 2242 WITTGENSTEIN, Ludwig. 1953. *Philosophical Investigations / Philosophische*
2243 *Untersuchungen*. Oxford: Basil Blackwell Publishers. Edited by G.E.M. Anscombe
2244 and R. Rhees, translation from the German by G.E.M. Anscombe; third edition:
2245 Wittgenstein (1967).
- 2246 —. 1967. *Philosophical Investigations / Philosophische Untersuchungen*. 3rd ed.
2247 Oxford: Basil Blackwell Publishers. Third edition of Wittgenstein (1953),
2248 edited by G.E.M. Anscombe and R. Rhees, translation from the German by
2249 G.E.M. Anscombe.
- 2250 ZAGZEBSKI, Linda Trinkaus. 1996. *Virtues of the Mind: An Inquiry into the Nature of*
2251 *Virtue and the Ethical Foundations of Knowledge*. Cambridge: Cambridge Univer-
2252 sity Press, doi:10.1017/cbo9781139174763.

The Dual-Detector Argument Against the Modal Theory

DAN MARSHALL

The modal theory holds that facts (properties) are identical iff they are necessarily equivalent (coextensive). One of the most prominent arguments against the modal theory is Elliot Sober's dual-detector argument. According to this argument, the fact that some particular thing is a triangle is distinct from the necessarily equivalent fact that it is a trilateral, since it is only the former fact that causes an output of a certain machine. I argue that the dual-detector argument fails, in part because whatever initial plausibility it has relies on the failure to take into consideration a needed relativisation to times and the failure to distinguish between two facts collectively causing a fact and their conjunction singly causing it. I also argue that variants of the argument are equally unsuccessful.

One of the most popular and well known accounts of the identity-conditions of facts and properties is the modal theory.¹ According to this theory: i) two facts are identical iff they are necessarily equivalent to each other; and ii) two properties are identical iff they are necessarily coextensive to each other. That is, the modal theory holds that: i) the fact that ϕ = the fact that ψ iff, necessarily, (ϕ iff ψ); and ii) the property of being F = the property of being G iff, necessarily, for any x , (x is F iff x is G).² This theory is prima facie attractive, since it is simple to formulate and provides an account of the

1 A fact, as I will understand it here, is an obtaining state of affairs, where: i) a state of affairs is either a way things are or a way things aren't, and ii) a state of affairs obtains iff it is a way things are. A fact on this understanding is therefore just as way things are. Proponents of the modal theory include Stalnaker (1984), Lewis (1986a), and Jackson (1998).

2 For simplicity, I will assume necessitism, according to which, necessarily, for any x , necessarily, for some y , $x = y$. (Without this assumption, 'necessarily, for any x ' needs to be replaced with 'necessarily, for any x , necessarily' in the above characterisation of the modal theory.) I will also assume an abundant theory of facts and properties according to which (except for restrictions needed to avoid paradox), all true sentences express facts, and all predicates that can be used to form true or false sentences express properties.

2274 identity-conditions of facts and properties in terms of (at least relatively)
 2275 well understood notions. Everything else being equal, the theory is also more
 2276 parsimonious than rival theories that reject it, since, everything else being
 2277 equal, there are less facts if the modal theory holds than if it fails to hold and
 2278 there are distinct facts that are necessarily equivalent to each other.

2279 A prominent argument against the modal theory is the dual-detector ar-
 2280 gument originally due to Elliot Sober.³ Briefly, according to this argument,
 2281 there could be a machine that, as a result of containing detectors measuring
 2282 different aspects of an input, is causally sensitive to one fact without being
 2283 causally sensitive to another necessarily equivalent fact. Since, by Leibniz's
 2284 law, it follows from this that, contra the modal theory, there are distinct facts
 2285 that are necessarily equivalent to each other, the argument concludes that the
 2286 modal theory is false. Despite this argument's prominence, discussions of the
 2287 argument by both its proponents and opponents have been brief and cursory.
 2288 This paper will provide a more sustained evaluation of the dual-detector ar-
 2289 gument and will argue that such an evaluation shows that the argument is
 2290 unsuccessful.⁴

2291 I will proceed as follows. In section 1, I will formulate the dual-detector
 2292 argument before then arguing in section 2 that it is unsuccessful. In section 3,
 2293 I will then consider two variants of the argument and I will argue that these
 2294 variants, and more generally that all variants, are also unsuccessful.

2295 Before proceeding to section 1, it will be useful to briefly discuss another
 2296 common argument against the modal theory—the constituency argument—
 2297 in order to set it aside (see, for example, Audi 2016). Suppose (1) and (2) are
 2298 true, where 'W' refers to some particular wire.

- 2299 (1) *W* is a closed straight-sided figure that has three *angles*.
 2300 (2) *W* is a closed straight-sided figure that has three *sides*.

2301 According to the constituency argument, since the fact expressed by (1) has
 2302 angularity as a constituent while the fact expressed by (2) doesn't have this
 2303 property as a constituent, the facts expressed by (1) and (2) are not identical

3 See Sober (1982). A similar argument is also given by Enç (1982). A recent proponent of the dual-detector argument, for example, is Audi (2016). Other philosophers who are sympathetic to the argument include Miller (1995, 859) and Molnar (2003, 66). Opponents of the argument include Jackson (1998, 125–126) and Armstrong (1997, 145–146).

4 Two other arguments against the modal theory that appeal to causation have been given by Achinstein (1974) and Perry (1989). Sober (1982, 84–85) gives what I take to be a convincing response to Achinstein's argument. For a response to Perry's argument, see Marshall (2021).

2304 to each other. Since the modal theory entails that the facts expressed by (1)
 2305 and (2) are identical to each other (since they are necessarily equivalent to
 2306 each other), the constituency argument concludes from this that the modal
 2307 theory is false.

2308 The constituency argument arguably begs the question against the modal
 2309 theory by in effect assuming the rival structured theory of facts. According to
 2310 this rival theory, facts are structured in the same kind of way that sentences
 2311 are structured. In particular, according to the structured theory, facts are built
 2312 up out of objects, properties, relations, operators and quantifiers in the same
 2313 way that sentences are built up out of names, predicates, operator expressions
 2314 and quantifier expressions.⁵ If the structured theory holds so that the facts
 2315 expressed by (1) and (2) are built up out of objects, properties, relations,
 2316 operators and quantifiers in the same kind of way that sentences are built up
 2317 out of names, predicates and other expressions, then it is plausible that the
 2318 fact expressed by (1) has angularity as a constituent while the fact expressed
 2319 by (2) doesn't have this constituent. This is much less plausible, however, if
 2320 the structured theory is false and facts aren't structured like sentences. For
 2321 example, if facts are instead structured like visual experiences or pictures, then,
 2322 since it is prima facie plausible to associate (1) and (2) with the same (type) of
 2323 visual experience or picture, it is prima facie plausible that (1) and (2) express
 2324 the same fact and hence prima facie plausible that the facts expressed by (1)
 2325 and (2) don't differ in what constituents they have. (This is because it is at
 2326 least prima facie plausible that any picture that represents W as being a closed
 2327 straight-sided figure that has three *angles* also represents W as being a closed
 2328 straight-sided figure that has three *sides*, and vice versa.) Since the argument
 2329 from constituency provides no reason to think that facts are structured in
 2330 the way that the structured theory holds that they are structured, rather than
 2331 some other way, the argument therefore fails to provide a good reason to think
 2332 that (1) and (2) express distinct facts and hence fails to provide a good reason
 2333 to reject the modal theory.

5 The structured theory can be formulated more precisely as a thesis endorsing schemas such as (PS) and (OS) (see, for example, [Dorr 2016, 58–59](#)).

PS. For any x and y , if the fact that x is F = the fact that y is G , then: i) $x = y$, and ii) the property of being F = the property of being G .

OS. If the fact that $\pi_1(\phi_1)$ = the fact that $\pi_2(\phi_2)$, then: i) the operator of π_1 = the operator of π_2 , and ii) the state of affairs of it being that ϕ_1 = the state of affairs of it being that ϕ_2 .

2334 It is important to appreciate that the structured theory is neither self-evident
 2335 nor *prima facie* highly plausible, and hence it cannot simply be assumed to
 2336 hold in the above argument from constituency without begging the question
 2337 against the modal theory. Three brief reasons for this are the following: First,
 2338 prior to investigation and argument, the claim that facts are structured like
 2339 sentences is no more plausible than the claim that facts have some other type
 2340 of structure, such as that of visual experiences or pictures. Second, while (1)
 2341 and (2) arguably differ in their cognitive significance, since a linguistically
 2342 competent person arguably might endorse one of them while rejecting the
 2343 other, such a difference in cognitive significance is widely thought to be able to
 2344 be explained by a difference in what mode of presentation the facts expressed
 2345 by (1) and (2) have when expressed by these sentences, where this explanation
 2346 does not require that the facts expressed by these sentences are non-identical
 2347 (see, for example, Braun 1998; McKay and Nelson 2010). Third, the structured
 2348 theory conflicts with claims that are widely thought to be at least as *prima*
 2349 *facie* plausible as the structured theory itself, such as the claim made by (3).

2350 (3) ‘*W* is self-identical’ expresses the same fact as ‘*W* is identical to *W*’.

2351 (3) conflicts with the structured theory, since, if the structured theory is
 2352 true, the fact expressed by ‘*W* is self-identical’ has the property of being self-
 2353 identical as a constituent while the fact expressed by ‘*W* is identical to *W*’
 2354 lacks this constituent and instead has the property of being identical to *W* as a
 2355 constituent.⁶ Due to the above difficulty with the constituency argument, and
 2356 since we cannot simply assume the structured theory in arguing against the
 2357 modal theory, I will assume in the following that the constituency argument
 2358 against the modal theory fails.

2359 1 The Dual-Detector Argument

2360 The dual-detector argument is not meant to rely on the cogency of the con-
 2361 stituency argument discussed above, nor is it meant to rely on the truth of
 2362 the structured theory of facts. Instead, the dual-detector argument is meant
 2363 to provide a separate reason for rejecting the modal theory. The argument

6 A further possible consideration against the structured theory is that, unlike the modal theory, it gives rise to the Myhill-Russell paradox. See, for example, Dorr (2016) and Goodman (2017). For attempted solutions to the Myhill-Russell paradox that are compatible with the structured theory, see, for example, Walsh (2016), Kment (2022) and Yu (2017). See Dorr (2016) and Bjerring and Schwarz (2017) for further arguments against the structured theory.

2364 involves a machine *M* that contains two detectors: a closed straight-sided
2365 figure detector and a three-angle detector. These detectors are linked in a
2366 series in *M*, so that, if a wire (or several wires) are inputted into *M*, they are
2367 first inputted into the closed straight-sided figure detector and then, if they are
2368 outputted by this first detector, they are inputted into the three-angle detector.
2369 If the wire (or wires) are then outputted by the three-angle detector, they are
2370 then outputted by *M*. Indeed, I will assume in the following that what it is for
2371 something (or some things) to be outputted by *M* is just for it (or them) to be
2372 outputted by this second detector.

2373 The closed straight-sided figure detector in *M* works so that “when given a
2374 piece of wire as input, it will output the piece of wire if and only if the wire is
2375 a closed [(plane)] figure and all sides of the figure are straight” (Sober 1982,
2376 185). More explicitly, let us say that: i) when given a piece of wire as input
2377 that is a closed figure all of whose sides are straight, the closed straight-sided
2378 figure detector outputs the wire, and it does this *because* the wire is a closed
2379 figure all of whose sides are straight; whereas, ii) when given a piece of wire
2380 (or several pieces of wire) as input that is not a single closed figure all of
2381 whose sides are straight, the closed-straight-sided figure detector does not
2382 output it (or them). The three-angle detector, on the other hand, works so
2383 that “when given any number of straight pieces of wire, it outputs them if
2384 and only if they have three angles” (Sober 1982, 185). More explicitly: i) when
2385 given one or more pieces of wire with straight sides that collectively have
2386 three angles, the three-angle detector outputs them and it does this *because*
2387 the wire (or wires) collectively have three angles; whereas, ii) when given
2388 one or more pieces of wire with straight sides that don’t collectively have
2389 three angles, the three-angle detector does not output them. The three-angle
2390 detector is causally sensitive to whether the input has three angles, and not
2391 to whether it has three sides, since, when given a four-sided open figure, it
2392 will output the object (since it has three angles), and it will fail to do this if
2393 the four-sided figure is closed. In addition, when the three-angle detector is
2394 given three unconnected pieces of wire, each containing exactly one angle,
2395 the detector will output them, even though it is made up of six straight line
2396 segments.

2397 Sober states the dual-detector argument as follows:

2398 Now consider a particular object—a piece of wire—which is fed
2399 into the machine, passes through both [detectors], and is then
2400 outputted by the machine. What property of the object *caused* it

2401 to be outputted? Given the mechanism at work here, I think that
 2402 the cause was the object's having the property of being a *closed*
 2403 *straight-sided figure having three angles* (i.e., its being a triangle),
 2404 and not its being a *closed straight-sided figure having three sides*
 2405 (i.e., its being a trilateral). If this is right, and if a difference in
 2406 causal efficacy is enough to insure a difference in property, it
 2407 follows that being a triangle is not the same property as being
 2408 a trilateral, even though "triangle" and "trilateral" are logically
 2409 (mathematically) equivalent. (Sober 1982, 185, author's emphasis)

2410 Let [ϕ] abbreviate 'the fact that ϕ ', and suppose that W is the piece of wire
 2411 that is fed into M . Let us also suppose that the above process of W being fed
 2412 into and then being sequentially outputted by the two detectors has occurred.
 2413 Then, according to Sober's dual-detector argument, (ANGLE) is true while
 2414 (SIDE) is false.

2415 ANGLE. [W is a closed straight-sided figure that has three *angles*]
 2416 causes [M outputs W].

2417 SIDE. [W is a closed straight-sided figure that has three *sides*] causes
 2418 [M outputs W].

2419 The dual-detector argument then employs Leibniz's law to infer from this
 2420 that, since they differ in what they cause, [W is a closed straight-sided figure
 2421 having three *angles*] is not identical to the necessarily equivalent [W is a
 2422 closed straight-sided figure having three *sides*]. The argument then infers
 2423 from (4) and the non-identity of these facts that the property of being a closed
 2424 straight-sided figure that has three *angles* (or being triangular) is not identical
 2425 to the necessarily coextensive property of being a closed straight-sided figure
 2426 that has three *sides* (or being trilateral).

2427 (4) For any x , IF x is F , x is G , and the property of being F = the property
 2428 of being G , THEN [x is F] = [x is G].

2429 Since these facts and properties are respectively necessarily equivalent to each
 2430 other and necessarily coextensive with each other (and hence are identical to

2431 each other according to the modal theory), the dual-detector argument then
 2432 concludes from the above results that the modal theory is false.⁷

2432 **Against the Dual-Detector Argument**

2434 One initial problem with the dual-detector argument is that (ANGLE) is not
 2435 strictly speaking true, at least if we assume as we did above that the above
 2436 described process involving *W* and *M* has already occurred.

2437 ANGLE. [*W* is a closed straight-sided figure that has three *angles*]
 2438 causes [*M* outputs *W*].

2439 To see why this is the case, let us suppose that, after being fed into *M* and put
 2440 inside the closed straight-sided figure detector at t_1 , *W* is outputted by the
 2441 closed straight-sided figure detector so that, at t_2 , *W* is inside the three-angle
 2442 detector. Let us also suppose that *W* being inside the three-angle detector at
 2443 t_2 results in *W* being outputted by the three-angle detector at t_3 , and hence
 2444 results in *W* being outputted by *M* at t_3 . Finally, let us also suppose that the
 2445 times t_1 , t_2 and t_3 are all past times. Then the fact that *W* is a closed straight-
 2446 sided figure that has three angles (either simpliciter or at the present time)
 2447 does not cause *M* to do anything to *W*, since *W* is no longer interacting with
 2448 *M*.

7 I am assuming that facts can cause other facts, where this claim is compatible with it also being the case that events can cause other events. If it is instead held that it is only events that can be causal relata, then 'fact' can be replaced with 'event' in the above presentation of the dual-detector argument to get the conclusion that there are distinct necessarily equivalent events (where two events are necessarily equivalent iff, necessarily, they either both occur or they both fail to occur). This conclusion together with (A) entails that there are distinct necessarily equivalent properties which, given (MF), entails that there are distinct necessarily equivalent states of affairs.

A. If the property of being *F* = the property of being *G*, then, for any *x*, the event of *x* having *F* = the event of *x* having *G*.

MF. The property of being *F* = the property of being *G* iff, necessarily, for any *x*, the state of affairs of *x* being *F* = the state of affairs of *x* being *G*.

Taking facts to be obtaining states of affairs (as in footnote 1), it follows from this that there are distinct necessarily equivalent facts.

2449 The above problem with the dual-detector argument shows that, as it is
 2450 most charitably understood, it is not (**ANGLE**) that is true according to the
 2451 argument, but is instead either (**ANGLE**_{*t*₁}) or (**ANGLE**_{*t*₂}).⁸

2452 **ANGLE**_{*t*₁}. [*At t*₁, *W* is a closed straight-sided figure that has three
 2453 *angles*] causes [*M* outputs *W* at *t*₃].

2454 **ANGLE**_{*t*₂}. [*At t*₂, *W* is a closed straight-sided figure that has three
 2455 *angles*] causes [*M* outputs *W* at *t*₃].

2456 As a result of this need to relativise to either time *t*₁ or time *t*₂, we therefore
 2457 have two versions of the dual-detector argument. The first version—the *t*₁-
 2458 version—holds that (**ANGLE**_{*t*₁}) is true and (**SIDE**_{*t*₁}) is false, from which it
 2459 infers that, contra the modal theory, the necessarily equivalent facts [*at t*₁, *W*
 2460 is a closed straight-sided figure that has three *angles*] and [*at t*₁, *W* is a closed
 2461 straight-sided figure that has three *sides*] are non-identical.

2462 **SIDE**_{*t*₁}. [*At t*₁, *W* is a closed straight-sided figure that has three *sides*]
 2463 causes [*M* outputs *W* at *t*₃].

2464 The second version of the dual-detector argument—the *t*₂-version—holds
 2465 instead that (**ANGLE**_{*t*₂}) is true and (**SIDE**_{*t*₂}) is false, from which it infers that,
 2466 contra the modal theory, the necessarily equivalent facts [*at t*₂, *W* is a closed
 2467 straight-sided figure that has three *angles*] and [*at t*₂, *W* is a closed straight-
 2468 sided figure that has three *sides*] are non-identical.

8 In response to the above problem with Sober's (1982) original formulation of the dual-detector argument, we might modify *M* so that its two detectors act on *W* at the same time rather than sequentially. Such a modified version of the argument faces the same difficulties as the *t*₁-version of the argument discussed below. First, given this modification, while it is plausible that (**ANGLE**^{*}) is true and (**SIDE**^{*}) is false (when relativised uniformly to the relevant time), there is an exclusion argument that argues from the truth of (**ANGLE**^{*}) to the falsity of (**ANGLE**).

ANGLE^{*}. [*W* is a closed straight-sided figure] and [*W* has three *angles*] collectively cause [*M* outputs *W*].

SIDE^{*}. [*W* is a closed straight-sided figure] and [*W* has three *sides*] collectively cause [*M* outputs *W*].

Second, this modified version of the argument faces the problem that, even if this exclusion argument is rejected, it doesn't seem possible to justify both the truth of (**ANGLE**) and the falsity of (**SIDE**).

2469 SIDE_{t₂}. [At t₂, W is a closed straight-sided figure that has three sides]
 2470 causes [M outputs W at t₃].

2471 As we will see, both these versions of the dual-detector argument have serious
 2472 problems.⁹

2473 The t₂-version of the dual-detector argument can be quickly seen to fail
 2474 as follows: It is [W has three angles at t₂] that causes W to be outputted by
 2475 the three-angle detector at t₃, rather than say [at t₂, W has three angles and
 2476 is blue] that causes this fact (even supposing that W is blue at t₂). This is
 2477 intuitively because [at t₂, W has three angles and is blue] goes beyond what
 2478 is causally relevant to whether W is outputted by the three-angle detector at
 2479 t₃. Similarly, it is [W has three angles at t₂] that causes W to be outputted by
 2480 the three-angle detector at t₃ rather than [at t₂, W is a closed straight-sided
 2481 figure that has three angles] that causes this fact. This is because the latter
 2482 fact also goes beyond what is causally relevant to whether W gets outputted
 2483 by the three-angle detector at t₃. Since W getting outputted by the three-angle
 2484 detector just is what it is for M to be outputted by W, it follows that (ANGLE_{t₂})
 2485 is false.

2486 ANGLE_{t₂}. [At t₂, W is a closed straight-sided figure that has three
 2487 angles] causes [M outputs W at t₃].

2488 Since the falsity of (ANGLE_{t₂}) conflicts with the t₂-version of the dual-detector
 2489 argument, this version of the argument fails.

2490 I will now argue that the t₁-version of dual-detector argument is also unsuccess-
 2491 ful and hence that both versions of the dual-detector argument fail. I will
 2492 do this by first giving an argument from causal exclusion that, contrary to the

9 There is also a temporally mixed version of the dual-detector argument that holds that (ANGLE_{t₁,t₂}) is true and (SIDE_{t₁,t₂}) is false.

ANGLE_{t₁,t₂}. [W is a closed straight-sided figure at t₁ that has three angles at t₂] causes [M outputs W at t₃].

SIDE_{t₁,t₂}. [W is a closed straight-sided figure at t₁ that has three sides at t₂] causes [M outputs W at t₃].

This version of the argument at best only shows that [W is a closed straight-sided figure at t₁ that has three angles at t₂] is not identical to [W is a closed straight-sided figure at t₁ that has three sides at t₂], which does not conflict with the modal theory since these facts are not necessarily equivalent to each other.

2493 dual-detector argument, (ANGLE_{t_1}) is false. I will then argue that, even if this
 2494 causal exclusion argument is rejected, it is not possible to justify both the truth
 2495 of (ANGLE_{t_1}) and the falsity of (SIDE_{t_1}), the justification of both of which is
 2496 required for the t_1 -version of the argument to be successful. (Or at least, I will
 2497 argue that one cannot justify the truth of (ANGLE_{t_1}) and the falsity of (SIDE_{t_1})
 2498 without appealing to some other argument against the modal theory that, if
 2499 successful, would refute the modal theory by itself and hence would render
 2500 the dual-detector argument superfluous.)

2501 To set up the needed background for the argument from causal exclusion
 2502 against (ANGLE_{t_1}), note that, in the case of M processing W , [W is a closed
 2503 straight-sided figure at t_1] causes W to be outputted by the closed straight-
 2504 sided figure detector, and so causes W to be in the three-angle detector at t_2 .
 2505 Hence we have (5)

2506 (5) [W is a closed straight-sided figure at t_1] causes [W is in the three-angle
 2507 detector at t_2].

2508 Since [W is in the three-angle detector at t_2] and [W has three angles at t_2]
 2509 collectively cause W to be outputted by the three-angle detector at t_3 , which
 2510 is what it is to be outputted by M at t_3 , we also have (6).

2511 (6) [W is in the three-angle detector at t_2] and [W has three angles at t_2]
 2512 collectively cause [M outputs W at t_3].

2513 Since plausibly one of the causes of W having three angles at t_2 is that it had
 2514 three angles at previous times before t_2 , (7) plausibly also holds.

2515 (7) [W has three angles at t_1] causes [W has three angles at t_2].

2516 Assuming, as is plausible, that the causal transitivity principle (T) holds in
 2517 this causal situation, (5–7) then entail ($\text{ANGLE}_{t_1}^*$).¹⁰

2518 T. IF the members of Φ_1 collectively cause r_1 , the members of Φ_2
 2519 collectively cause $r_2 \dots$ and $r_1, r_2 \dots$ collectively cause r ; THEN the
 2520 members of $\Phi_1 \cup \Phi_2 \cup \dots$ collectively cause r .

10 While causation is plausibly transitive in many typical cases, such as in the case above, many philosophers hold that causation is not unrestrictedly transitive. For alleged counterexamples to transitivity, see, for example, Kwart (1991) and McDermott (1995). For a defense of transitivity unrestrictedly holding, see Hall (2000).

2521 ANGLE_{t₁}*. [*W* is a closed straight-sided figure at *t*₁] and [*W* has three
2522 angles at *t*₁] collectively cause [*M* outputs *W* at *t*₃].

2523 With the above background in place, it might seem like it should now be
2524 easy to derive (ANGLE_{t₁}) from (ANGLE_{t₁}*), and hence establish that (ANGLE_{t₁})
2525 holds.

2526 ANGLE_{t₁}. [At *t*₁, *W* is a closed straight-sided figure that has three
2527 angles] causes [*M* outputs *W* at *t*₃].

2528 However using the above background, we can now give the following argu-
2529 ment from causal exclusion that (ANGLE_{t₁}) is instead false: Just as [at *t*₂, *W* is
2530 a closed straight-sided figure that has three angles] fails to cause the closed
2531 straight-sided figure detector to output *W* at *t*₃ (since the former fact goes
2532 beyond what is causally relevant), [at *t*₁, *W* is a closed straight-sided figure
2533 that has three angles] fails to cause the closed straight-sided figure detector
2534 to output *W* (since this fact also goes beyond what is causally relevant) and
2535 hence this fact fails to cause [*W* is in the three-angle detector at *t*₂]. Hence we
2536 have (8).

2537 (8) [at *t*₁, *W* is a closed straight-sided figure that has three angles] does not
2538 cause [*W* is in the three angle detector at *t*₂].

2539 Similarly, while [*W* has three angles at *t*₁] is a cause of [*W* has three angles
2540 at *t*₂], it is not the case that [at *t*₁, *W* is a closed straight-sided figure that has
2541 three angles] causes this fact, since it goes beyond what is causally relevant.
2542 Hence we have (9).

2543 (9) [at *t*₁, *W* is a closed straight-sided figure that has three angles] does not
2544 cause [*W* has three angles at *t*₂].

2545 Since [at *t*₁, *W* is a closed straight-sided figure that has three angles] is also
2546 not caused by either [*W* is in the three-angle detector at *t*₂] or [*W* has three
2547 angles at *t*₂], and the causal chain that leads up to [*M* outputs *W* at *t*₃] goes
2548 through [*W* is in the three-angle detector at *t*₂] and [*W* has three angles at *t*₂],
2549 it therefore follows from (8) and (9) that [at *t*₁, *W* is a closed straight-sided
2550 figure that has three angles] isn't part of the causal chain that leads to [*M*
2551 outputs *W* at *t*₃] and hence does not cause it. Hence (ANGLE_{t₁}) is false.

2552 ANGLE_{t₁}. [At t₁, W is a closed straight-sided figure that has three
2553 angles] causes [M outputs W at t₃].

2554 A more rigorous version of the above argument against (ANGLE_{t₁}) can be
2555 given by appealing to the version of the principle of causal exclusion given by
2556 (PCE).¹¹

2557 PCE. In cases where there is no genuine causal overdetermination,
2558 if S is a set of facts that occur at a time t whose members collectively
2559 completely cause f, then S is the unique set of facts that occur at t
2560 and collectively completely cause f.

2561 In (PCE), a fact is said to occur at a certain time iff the fact only concerns
2562 how things are at that time. Genuine causal overdeterminism, on the other
2563 hand, occurs when two independent causal processes converge on the same
2564 effect, such as when a house burns down because a lit match starts a fire in
2565 the garbage at the same time as lightning strikes the house.

2566 Since there is no genuine causal overdetermination in the case of W being
2567 outputted by M, (PCE) can be used to argue that (ANGLE_{t₁}) is false as follows:
2568 Suppose, for reductio, that (ANGLE_{t₁}) is true. Then [W is a closed straight-
2569 sided figure that has three angles at t₁] together with the members of some
2570 possibly empty set Ψ₁ completely cause [M outputs W at t₃]. Since (ANGLE_{t₁}^{*})
2571 holds, it is also true that [W is a closed straight-sided figure at t₁], [W has
2572 three angles at t₁] together with the members of some possibly empty set Ψ₂
2573 collectively completely cause [M outputs W at t₃]. Since the relevant facts
2574 occur at the same time, these two consequences together with (PCE) then
2575 entail (10).

2576 (10) [W is a closed straight-sided figure at t₁], [W has three angles at t₁], [W
2577 is a closed straight-sided figure that has three angles] and the members
2578 of some possibly empty set Ψ collectively completely cause [M outputs
2579 W at t₃].

2580 If (10) is true, then [W is a closed straight-sided figure at t₁], [W has three
2581 angles at t₁] and the members of Ψ by themselves collectively completely
2582 cause [M outputs W at t₃], since [at t₁, W is closed straight-sided figure that

11 For discussion of the principle of causal exclusion, see, for example, Kim (2005) and Moore (2018).

2583 has three angles] is superfluous given the presence of [W is a closed straight-
 2584 sided figure at t_1] and [W has three angles at t_1]. Given (PCE), however, this
 2585 consequence conflicts with (10). Hence, the reductio assumption (ANGLE $_{t_1}$)
 2586 is false.

2587 The above argument shows that (ANGLE $_{t_1}$) fails to hold if (PCE) holds. Not
 2588 all philosophers, however, accept (PCE), and these philosophers will not be
 2589 convinced by the above argument from causal exclusion that the dual-detector
 2590 argument fails. For example, some philosophers reject (PCE) on the grounds
 2591 that it conflicts with the popular counterfactual dependency thesis (DEP).¹²

2592 DEP. Suppose that f and g obtain, and that, had f failed to obtain,
 2593 it would have been that g failed to obtain. Then, f causes g .

2594 Other philosophers reject (PCE) because they hold that, in cases where there
 2595 is no genuine causal overdetermination of a fact, there can still be multiple
 2596 complete causal chains that converge on that fact, provided these chains
 2597 are systematically related to each other in the right way. In particular, some
 2598 philosophers hold that there can be multiple such causal chains provided that,
 2599 for each such chain, either that chain generates all the other chains, or that
 2600 chain is generated by at least one other such chain. Someone who endorses
 2601 this view, for example, might endorse (CONJ).¹³

2602 CONJ. If f_1 and f_2 together with the members of a set Φ collectively
 2603 completely cause f , then the conjunction of f_1 and f_2 together with
 2604 the members of Φ collectively completely cause f .

2605 It follows from (CONJ) that, contra (PCE), if there is one causal chain leading
 2606 to f that contains the facts f_1 and f_2 occurring at a time t , then there is a further
 2607 causal chain which is systematically related to it by virtue of containing the
 2608 conjunction of f_1 and f_2 instead of f_1 and f_2 themselves. Given (CONJ), it is
 2609 natural to hold that this further causal chain containing the conjunction of f_1
 2610 and f_2 is generated by the former chain containing its conjuncts.

12 See, for example, Loewer (2007). Proponents of (DEP) typically place certain restrictions on (DEP), such as requiring that the counterfactual is to be read in a suitable non-backtracking sense (see Lewis 1973), that the facts (or events, when (DEP) is applied to events) that stand in the causation relation are “sufficiently distinct” (so that, for example, we don’t have the consequence that each fact causes itself) (see Kim 1973; Lewis 1986c), and that these facts (or events) are non-disjunctive (see Lewis 1986c).

13 Φ in (CONJ) can be the empty set.

2611 In light of the above views, the argument from causal exclusion does not by
 2612 itself decisively refute the t_1 -version of the dual-detector argument. In addition
 2613 to facing the argument from causal exclusion, however, the t_1 -version of the
 2614 dual-detector argument faces the problem that, even if the causal exclusion
 2615 argument fails, it doesn't appear possible to justify the truth of (ANGLE_{t_1})
 2616 while also justifying the falsity of (SIDE_{t_1}). (Or at least, it doesn't seem possible
 2617 to do this without relying on some other argument against the modal theory
 2618 which, if successful, would by itself refute the modal theory. I will discuss two
 2619 attempts to give such a justification, and I will argue that both these attempts
 2620 fail. The failure of these two attempts will give us reason to think that no such
 2621 justification is possible, and hence reason to think that, even if (PCE) and
 2622 the argument from causal exclusion fail, the t_1 -version of the dual-detector
 2623 argument is still unsuccessful.

2624 The first attempt to justify the truth of (ANGLE_{t_1}) (while also justifying
 2625 the falsehood of (SIDE_{t_1}) appeals to (CONJ) above. This first attempt accepts
 2626 ($\text{ANGLE}_{t_1}^*$) on the basis of the transitivity reasoning given for it above. It then
 2627 infers from ($\text{ANGLE}_{t_1}^*$) and (CONJ) that the conjunction of [W is a closed
 2628 straight-sided figure at t_1] and [W has three angles at t_1] collectively (partially)
 2629 cause M to output W at t_3 . Assuming (as I will from now on) that this con-
 2630 junction is the fact [at t_1 , W is a closed straight-sided figure that has three
 2631 angles], it follows from this that (ANGLE_{t_1}) is true.

2632 ANGLE_{t_1} . [At t_1 , W is a closed straight-sided figure that has three
 2633 *angles*] causes [M outputs W at t_3].

2634 Let us assume that the above justification of (ANGLE_{t_1}) is successful. The
 2635 question that now needs to be addressed is whether we can go on to justify
 2636 the falsehood of (SIDE_{t_1}).

2637 SIDE_{t_1} . [At t_1 , W is a closed straight-sided figure that has three *sides*]
 2638 causes [M outputs W at t_3].

2639 One argument that tries to justify the falsehood of (SIDE_{t_1}) is the following:
 2640 Unlike (ANGLE_{t_1}), (SIDE_{t_1}) cannot be generated from the causal facts given to
 2641 us in the description of M processing W given in the dual-detector argument
 2642 using causal generational principles such as (T) and (CONJ). As a result, the
 2643 truth of (SIDE_{t_1}) would require some additional primitive causal fact to hold
 2644 in the case of M processing W , which would be unparsimonious. Moreover,

2645 since any such additional primitive causal fact would only contingently hold,
 2646 the possibility of such a fact holding can be removed by simply stipulating
 2647 that no such additional primitive causal fact holds in the possible case of *M*
 2648 processing *W* that we are concerned with. Hence, according to this argument,
 2649 the truth of (*SIDE*_{*t*₁}) can be ruled out either on parsimony grounds or by
 2650 stipulation.

2651 The problem with this argument for the falsity of (*SIDE*_{*t*₁}) is that it begs the
 2652 question against the modal theory. It does this because, if the modal theory
 2653 is true, then, contra the above argument, (*SIDE*_{*t*₁}) can be generated from the
 2654 causal facts given to us in the description of the case of *M* processing *W* in the
 2655 dual-detector argument and the generational principles (*T*) and (*CONJ*) in the
 2656 same way that (*ANGLE*_{*t*₁}) can be so generated. This is because, if the modal
 2657 theory is true, then [at *t*₁, *W* is a closed straight-sided figure that has three
 2658 sides] is the conjunction of [*W* is a closed straight-sided figure at *t*₁] and [*W*
 2659 has three angles at *t*₁], just as much as [at *t*₁, *W* is a closed straight-sided figure
 2660 that has three angles] is. Hence, if the modal theory is true, then (*SIDE*_{*t*₁}) can
 2661 be derived from (*ANGLE*_{*t*₁}^{*}) and (*CONJ*) in the same way that (*ANGLE*_{*t*₁}) can.

2662 An alternative way of trying to justify the falsehood of (*SIDE*_{*t*₁}) appeals to
 2663 (*CONJ*^{*}).¹⁴

2664 *CONJ*^{*}. If the conjunction of *f*₁ and *f*₂ partially causes *f*, then *f*₁
 2665 and *f*₂ collectively partially cause *f*.

2666 We can give the same kind of argument from parsimony and contingency for
 2667 the falsity of (*SIDE*_{*t*₁}^{*}) as was given above for the falsity of (*SIDE*_{*t*₁}), with the
 2668 difference that this argument for the falsity of (*SIDE*_{*t*₁}^{*}), unlike the argument
 2669 for the falsity of (*SIDE*_{*t*₁}), does not beg the question against the modal theory.

14 (*CONJ*) and (*CONJ*^{*}) are in the vicinity of two principles, (*A*) and (*B*), that Sober appeals to when defending the dual-detector argument.

A. If two devices, “which are linked in series in the [machine], are sensitive just to properties *P* and *Q*, respectively, then the [machine] itself is sensitive to the conjunctive property *P*-and-*Q*.” (Sober 1982, 186)

B. If “two devices which are linked in series are such that the first is sensitive to *P* and the second is *not* sensitive to *R* (where *P* ≠ *R*, and neither implies the other), then the [machine] is *not* sensitive to the conjunctive property *P*-and-*R*.” (Sober 1982, 186)

As argued below in the case of (*CONJ*^{*}), (*B*) immediately conflicts with the modal theory and is hard to justify.

2670 $\text{SIDE}_{t_1}^*$. [W is a closed straight-sided figure at t_1] and [W has three
2671 sides at t_1] collectively cause [M outputs W at t_3].

2672 Indeed, plausibly both opponents and proponents of the modal theory should
2673 reject ($\text{SIDE}_{t_1}^*$). Given the falsity of ($\text{SIDE}_{t_1}^*$), however, the falsity of (SIDE_{t_1})
2674 follows from (CONJ^*).¹⁵ If we are justified in endorsing (CONJ^*), then, we can
2675 use it to justify the falsehood of (SIDE_{t_1}).

2676 One problem with (CONJ^*) is that the principle directly conflicts with the
2677 modal theory. This is because, if the modal theory holds, then (CONJ^*) has
2678 the absurd consequence that, if f partially causes g , then any fact h that
2679 is necessitated by f also causes g . (This is because, according to the modal
2680 theory, if a fact f necessitates a fact h , then f is the conjunction of f and
2681 h .) If [Suzy throws a rock] causes [the window breaks], for example, then,
2682 if the modal theory holds, (CONJ^*) entails that [Suzy throws a rock or Suzy
2683 does not throw a rock] (which is necessitated by [Suzy throws a rock]) also
2684 causes [the window breaks], which is absurd. In light of this, one problem
2685 with (CONJ^*) is that, if it is accepted, then we don't need the dual-detector
2686 argument to refute the modal theory, since (CONJ^*) by itself achieves this
2687 task. If the dual-detector argument needs to rely on (CONJ^*) in order to be
2688 successful, then, the argument is superfluous.

2689 A second (more serious) problem with (CONJ^*) is that it is not clear why
2690 we should believe it. A proponent of (CONJ^*) might attempt to justify the
2691 principle by arguing that, in ordinary language, sentences of the form (11)
2692 are equivalent to sentences of the form (12).

2693 (11) φ because ϕ and ψ .

2694 (12) φ because ϕ and because ψ .

2695 Such a proponent might then argue that (on its relevant causal use) (11) is
2696 equivalent to (11*) and (12) is equivalent to (12*).

2697 (11*) [ϕ and φ] causes [φ].

2698 (12*) [ϕ and φ] collectively cause [φ].

2699 Assuming that these equivalences all hold, it follows that (11*) entails (12*),
2700 from which it follows that (CONJ^*) holds.

15 I am assuming that [at t_1 , W is a closed straight-sided figure that has three sides] is the conjunction of [W is a closed straight-sided figure at t_1] and [W has three sides at t_1].

2701 A problem with this attempted justification for (CONJ*) is that (12) is plausi-
 2702 bly ambiguous between a conjunctive reading and a non-conjunctive reading,
 2703 just like (13) is (cf. Marshall 2021, 8035).

2704 (13) Jane wants to go swimming and go hiking.

2705 (13) has a non-conjunctive reading on which the proposition Jane is described
 2706 as desiring is the proposition that Jane goes swimming and hiking. On this
 2707 reading, (13) is true iff (13n) is true.

2708 (13n) Jane wants to go (swimming and hiking).

2709 (13) also has a conjunctive reading on which (13) is true iff (13c) is true.

2710 (13c) Jane wants to go swimming and Jane wants to go hiking.

2711 (11) is plausibly similarly ambiguous between a non-conjunctive reading
 2712 on which it is equivalent to (11n) and a conjunctive reading on which it is
 2713 equivalent to (11c).

2714 (11n) φ because (ϕ and ψ).

2715 (11c) (φ because ϕ) and (φ because ψ).

2716 On its conjunctive reading, while (11) is equivalent to ((12) (on its causal use),
 2717 there is no reason to think that (on its causal use) (11) is equivalent to (11*)
 2718 (or at least no such reason has yet been provided).¹⁶ On its non-conjunctive
 2719 reading, on the other hand, there is no reason to think that (11) is equivalent
 2720 to (12). As a result, appealing to natural language does not appear to help a
 2721 proponent of the dual-detector argument justify (CONJ*). In light of this, it
 2722 is not clear how (CONJ*) might be justified.¹⁷ As a result, it does not appear
 2723 possible to justify the truth of (ANGLE_{t₁}) by appealing to (CONJ) while also
 2724 justifying the falsity of (SIDE_{t₁}).

16 The claim that (12) is equivalent to (12*) can also be resisted, since it might be denied that 'f causes h' and 'g causes h' entails 'f and g collectively cause h'. For example, this inference might be thought to fail if f and g are individually complete causes of h that concern different times.

17 Or at least, it is not clear how (CONJ*) might be justified without begging the question against the modal theory. It might perhaps be possible to justify (CONJ*) if we assume the structured theory and give a general account of how less fundamental facts get to have their causal features in terms of the causal features of more fundamental facts that involves principles like (CONJ).

2725 I will discuss one further attempt to justify both the truth of (**ANGLE**_{*t*₁}) and
 2726 the falsity of (**SIDE**_{*t*₁}). Instead of appealing to (**CONJ**), this second attempt
 2727 appeals to the popular counterfactual dependency thesis (**DEP**) stated above.¹⁸

2728 DEP. Suppose that *f* and *g* obtain, and that, had *f* failed to obtain,
 2729 it would have been that *g* failed to obtain. Then, *f* causes *g*.

2730 Assuming that (**DEP**) holds, we can derive (**ANGLE**_{*t*₁}) as follows: In the case
 2731 of *M* outputting *W*, had it not been that, at *t*₁, *W* was a closed straight-sided
 2732 figure that had three *angles*, then either: i) *W* would not have been a closed
 2733 straight-sided figure at *t*₁; or ii) *W* would not have had three angles at *t*₁, in
 2734 which case *W* would also not have had three angles at *t*₂. If *W* had failed to be
 2735 a closed straight-sided figure at *t*₁, *W* would not have been outputted by the
 2736 closed straight-sided figure detector at *t*₂, and hence *W* would not have been
 2737 outputted by *M* at *t*₃. On the other hand, if *W* had failed to have three angles
 2738 at *t*₂, it would not have been outputted by the three-angle detector at *t*₃, and
 2739 hence would also not have been outputted by *M* at *t*₃. Hence, had it not been
 2740 that, at *t*₁, *W* was a closed straight-sided figure that had three *angles*, *M* would
 2741 not have outputted *W* at *t*₃. It therefore follows from (**DEP**) that (**ANGLE**_{*t*₁}) is
 2742 true.

2743 **ANGLE**_{*t*₁}. [*At t*₁, *W* is a closed straight-sided figure that has three
 2744 *angles*] causes [*M* outputs *W* at *t*₃].

2745 Assuming that (**DEP**) holds, then, a proponent of the dual-detector argument
 2746 can use (**DEP**) to justify (**ANGLE**_{*t*₁}). Unfortunately for proponents of the dual-
 2747 detector argument, however, if (**DEP**) holds it can also be used to justify the
 2748 truth of (**SIDE**_{*t*₁}). To see why, note that, had it not been that, at *t*₁, *W* was a
 2749 closed straight-sided figure that had three *sides*, then *W* would also either:
 2750 i) not have been a closed straight-sided figure at *t*₁ or ii) not have had three
 2751 angles at *t*₁, in which case it would not have had three angles at *t*₂. Hence, had
 2752 it not been that, at *t*₁, *W* was a closed straight-sided figure having three *sides*,
 2753 at least one of the detectors would not have outputted *W*, and so *M* would
 2754 not have outputted *W* at *t*₃. Hence, it also follows from (**DEP**) that (**SIDE**_{*t*₁})
 2755 is true. Hence, a proponent of the dual-detector argument cannot use (**DEP**)

18 Related principles we might try to appeal to in order to simultaneously justify the truth of (**ANGLE**) and the falsity of (**SIDE**) (which have similar problems to (**DEP**)) are difference-making principles, such as those proposed by Sartorio (2005) and List and Menzies (2009).

2756 to justify the combination of (ANGLE_{t_1}) being true and (SIDE_{t_1}) being false.
 2757 This second attempt at justifying the truth of (ANGLE_{t_1}) and the falsehood of
 2758 (SIDE_{t_1}) therefore fails.

2759 I have now discussed two attempts to justify the truth of (ANGLE_{t_1}) and the
 2760 falsity of (SIDE_{t_1}), and I have argued that both of these attempts fail. As far as
 2761 I can see, other attempts to do this are equally unsuccessful. If this is the case,
 2762 then both the t_1 -version and the t_2 -version of the dual-detector argument fail.

2763 Variants of the Dual-Detector Argument

2764 In the face of the failure of the original version of Sober's dual-detector argu-
 2765 ment, it might be thought that the argument can be modified so that it evades
 2766 the problems discussed in section 2. In particular, it might be thought that
 2767 these problems can be evaded by replacing the necessarily equivalent facts
 2768 expressed by ($1t_1$) and ($2t_1$) with some other necessarily equivalent facts and
 2769 describing a machine that is causally sensitive to one of these facts but not
 2770 the other.

2771 ($1t_1$) W is a closed straight-sided figure that has three *angles* at t_1 .

2772 ($2t_1$) W is a closed straight-sided figure that has three *sides* at t_1 .

2773 As far as I can see, however, this cannot be done.

2774 To illustrate the difficulty involved in successfully modifying the dual-
 2775 detector argument in the above manner, I will briefly consider two attempts
 2776 to do this that replace the facts expressed by ($1t_1$) and ($2t_1$) with the facts
 2777 expressed by (14) and (15), where W^* is a circular wire and where the facts
 2778 expressed by (14) and (15) are both necessarily equivalent to the fact that W^*
 2779 is a circle.¹⁹

2780 (14) W^* is a closed (plane) figure all of whose points are equidistant from a
 2781 point.

2782 (15) W^* is a closed (plane) figure of constant curvature.

2783 For the first attempt, consider a machine M_1^* that, when given a closed (plane)
 2784 figure as an input, scans that figure by having a distinct curvature detector for
 2785 each point of the figure. Suppose that each of these detectors measures the
 2786 curvature of their associated point in the figure and sends the result of this
 2787 measurement in the form of a signal to the CPU of M_1^* . Further, suppose that,

19 This variant was suggested by a referee.

2788 if all the signals the CPU receives are of the same value, then the fact that the
 2789 signals it receives have the same value causes the figure to be outputted by
 2790 M_1^* . Finally, suppose that the circular wire W^* is inputted into this machine
 2791 M_1^* , is scanned by it, and is then outputted by it. It might then be claimed that,
 2792 in this case, (CURV) is true while (DIST) is false, and that, due to Leibniz's law,
 2793 this difference in truth-value entails that the modal theory is false.

2794 CURV. [W^* is a closed figure with constant curvature] causes [M_1^*
 2795 outputs W^*].

2796 DIST. [W^* is a closed figure all of whose points are equidistant from
 2797 a point] causes [M_1^* outputs W^*].

2798 A problem with this first attempt at finding a successful variant of the dual-
 2799 detector argument is that it is no more obvious that (CURV) holds than it is
 2800 that (ANGLE $_{t_1}$) holds in Sober's original case.

2801 ANGLE $_{t_1}$. [At t_1 , W is a closed straight-sided figure that has three
 2802 angles] causes [M outputs W at t_3].

2803 Instead, using transitivity reasoning, what can be uncontroversially estab-
 2804 lished in the variant case of machine M_1^* is a claim along the lines of (CURV*),
 2805 just as what can be uncontroversially established using such reasoning in
 2806 Sober's original case of machine M is (ANGLE* $_{t_1}$).

2807 CURV*. [Point p_1 of W^* has curvature C], [point p_2 of W^* has
 2808 curvature C]... collectively cause [M_1^* outputs W^*].

2809 ANGLE* $_{t_1}$. [W is a closed straight-sided figure at t_1] and [W has three
 2810 angles at t_1] collectively cause [M outputs W at t_3].

2811 Moreover, an opponent of the modal theory who wishes to show that (CURV)
 2812 and (DIST) differ in their truth-value faces the same challenges that a pro-
 2813 ponent of Sober's original version of the dual-detector argument faces in
 2814 showing that (ANGLE $_{t_1}$) and (SIDE $_{t_1}$) differ in their truth-value. First, they
 2815 need to resist an argument from causal exclusion that (CURV*) entails the
 2816 falsehood of (CURV). And second, they need to find some way of justifying the
 2817 truth of (CURV) while also justifying the falsehood of (DIST), a task that ap-
 2818 pears to be just as difficult as finding a way of justifying the truth of (ANGLE $_{t_1}$)

2819 while also justifying the falsehood of (SIDE_{t_1}). Hence, this first attempt at
2820 describing a machine that is differentially sensitive to the facts expressed by
2821 (14) and (15) results in a variant of the dual-detector argument that is no more
2822 successful than Sober's original argument.

2823 For a second attempt to show that there could be a machine that is causally
2824 sensitive to one of the facts expressed by (14) and (15) but not the other,
2825 consider a machine M_2^* that contains an extendable straight rod that rotates
2826 around one of its endpoints. When given a closed figure as input, M_2^* works
2827 by placing this rod inside the inputted closed figure, fixing the location of
2828 one of the rod's endpoints, extending the length of the rod until its other
2829 endpoint touches the inputted figure, and then rotating the rod around its
2830 fixed endpoint while keeping the length of the rod fixed. If the rod does a full
2831 rotation without moving the inputted figure or losing touch with it, then the
2832 fact that it does this causes M_2^* to output the figure. Suppose now that the
2833 circular wire W^* is inputted into M_2^* and that the rod of M_2^* is placed inside
2834 of W^* and does a full rotation meeting the above conditions, so that W^* gets
2835 outputted by M_2^* . It might then be claimed that, in this case, (DIST) is true
2836 and (CURV) is false, and hence that the modal theory is false.

2837 The problem with this second variant of Sober's version of the dual-detector
2838 argument is that, if W^* is a circle that is inputted into and then outputted by
2839 M_2^* , then there is no reason to think that (CURV) and (DIST) differ in their
2840 truth-value. In particular, if W^* is so inputted and outputted, it is equally
2841 plausible to say that the machine measures the curvature of the points of W^*
2842 as it is to say that it measures the equidistance of those points from a common
2843 point. After all, the rod would fail to do its full rotation (while touching but
2844 not moving W^*) if the points of W^* didn't have constant curvature, just as it
2845 would fail to do this if the points of W^* weren't equally distant from some
2846 common point. There is therefore no grounds for thinking that W^* being
2847 outputted by M_2^* is due to one of these facts rather than the other. Hence, M_2^*
2848 also fails to be a demonstrable case of a machine that is causally sensitive to
2849 one of the facts expressed by (14) and (15) and not the other.

2850 Other variations of Sober's original version of the dual-detector argument
2851 face similar problems to those described above. Indeed, the above two attempts
2852 to construct a successful variant of Sober's original version of the argument
2853 arguably illustrate a dilemma facing any such attempt. This dilemma is the
2854 following: Suppose we have a machine whose output is intended to be caused
2855 by the fact f_1 and not by the necessarily equivalent fact f_2 . Then the machine
2856 will either contain multiple detectors that differ in what aspects of the input

2857 they measure (as in the cases of M and M_1^*), or the machine will only contain
 2858 detectors (or a single detector) that don't so differ (as in the case of M_2^*). If
 2859 the machine contains multiple detectors that differ in what aspects of the
 2860 input they measure, then the argument against the modal theory based on
 2861 this machine will arguably face the same challenges facing Sober's original arg-
 2862 gument and the first variant of it discussed above. In particular, the argument
 2863 will need to resist an argument from causal exclusion and will face the same
 2864 difficulties in justifying the claim that the input being outputted is caused by
 2865 f_1 and not by f_2 that Sober's original dual-detector argument faces in justifying
 2866 the truth of ($ANGLE_{t_1}$) and the falsity of ($SIDE_{t_1}$). On the other hand, if the
 2867 machine contains only a single detector (or multiple detectors that don't differ
 2868 in what aspects of the input they measure), then it will arguably fail to be
 2869 even initially plausible that f_1 and f_2 differ in whether they cause the input
 2870 to be outputted just as there is no even initial plausibility for thinking that
 2871 the facts expressed by (14) and (15) differ in whether they cause W^* to be
 2872 outputted by machine M_2^* . Hence, whether or not we have a machine that
 2873 contains detectors that differ in what aspects they measure, the argument
 2874 against the modal theory based on this machine will arguably fail. In light of
 2875 this, it is reasonable to conclude that, not only does Sober's original version of
 2876 the dual-detector argument fail, but it is not possible to modify the argument
 2877 so that it is successful. If this is correct, then all variants of the dual-detector
 2878 argument fail and some other kind of argument will be needed if we are to
 2879 have reason to reject the modal theory of facts and properties.*

2880 Dan Marshall

2881  0000-0002-5763-3875

2882 Lingnan University

2883 danmarshall@ln.edu.hk

2884 References

- 2885 ACHINSTEIN, Peter. 1974. "The Identity of Properties." *American Philosophical Quar-*
 2886 *terly* 11(4): 257–275.
- 2887 ARMSTRONG, David M. 1997. *A World of States of Affairs*. Cambridge: Cambridge
 2888 University Press, doi:10.1017/cbo9780511583308.

* Research in this paper was supported by an Early Career Scheme grant from the Research Grants Council of Hong Kong SAR, China (LU23607616). Thanks to Andrew Brenner, Daniel Waxman and three anonymous referees for their valuable comments on this paper.

- 2889 AUDI, Paul. 2016. "Property Identity." *Philosophy Compass* 11(12): 829–840, doi:10.1111
2890 1/phc3.12380.
- 2891 BJERRING, Jens Christian and SCHWARZ, Wolfgang. 2017. "Granularity Problems."
2892 *The Philosophical Quarterly* 67(266): 22–37, doi:10.1093/pq/pqw028.
- 2893 BRAUN, David. 1998. "Understanding Belief Reports." *The Philosophical Review* 107(4):
2894 555–595, doi:10.2307/2998375.
- 2895 DORR, Cian. 2016. "To Be *F* Is to Be *G*." in *Philosophical Perspectives 30: Metaphysics*,
2896 edited by John HAWTHORNE and Jason TURNER, pp. 39–134. Hoboken, New Jersey:
2897 John Wiley and Sons, Inc., doi:10.1111/phpe.12079.
- 2898 ENÇ, Berent. 2013. "Intentional States of Mechanical Devices." *Mind* 91(362): 161–183,
2899 doi:10.1093/mind/xci.362.161.
- 2900 GOODMAN, Jeremy. 2017. "Reality Is Not Structured." *Analysis* 77(1): 43–53, doi:10.1
2901 093/analys/anw002.
- 2902 HALL, Ned. 2000. "Causation and the Price of Transitivity." *The Journal of Philosophy*
2903 97(4): 198–222, doi:10.2307/2678390.
- 2904 JACKSON, Frank. 1998. *From Metaphysics to Ethics: A Defense of Conceptual Analysis*.
2905 Oxford: Oxford University Press, doi:10.1093/0198250614.001.0001.
- 2906 KIM, Jaegwon. 1973. "Causes and Counterfactuals." *The Journal of Philosophy* 70(17):
2907 570–572, doi:10.2307/2025312.
- 2908 —. 2005. *Physicalism, or Something Near Enough*. Princeton Monographs in Philosophy
2909 n. 19. Princeton, New Jersey: Princeton University Press, doi:10.1515/9781400840
2910 847.
- 2911 KMENT, Boris. 2022. "Russell-Myhill and Grounding." *Analysis* 82(1): 49–60, doi:10.1
2912 093/analys/anab028.
- 2913 KVART, Igal. 1991. "Transitivity and Preemption of Causal Relevance." *Philosophical*
2914 *Studies* 64(2): 125–160, doi:10.1007/bf00404826.
- 2915 LEWIS, David. 1973. "Causation." *The Journal of Philosophy* 70(17): 556–567. Reprinted,
2916 with a postscript (Lewis 1986d), in Lewis (1986b, 159–172), doi:10.2307/2025310.
- 2917 —. 1986a. *On the Plurality of Worlds*. Oxford: Blackwell Publishers.
- 2918 —. 1986b. *Philosophical Papers, Volume 2*. Oxford: Oxford University Press, doi:10.109
2919 3/0195036468.001.0001.
- 2920 —. 1986c. "Events." in *Philosophical Papers, Volume 2*, pp. 241–270. Oxford: Oxford
2921 University Press, doi:10.1093/0195036468.003.0008.
- 2922 —. 1986d. "Postscript to Lewis (1973)." in *Philosophical Papers, Volume 2*, pp. 172–213.
2923 Oxford: Oxford University Press, doi:10.1093/0195036468.003.0006.
- 2924 LIST, Christian and MENZIES, Peter. 2009. "Nonreductive Physicalism and the Limits
2925 of the Exclusion Principle." *The Journal of Philosophy* 106(9): 475–502, doi:10.584
2926 0/jphil2009106936.
- 2927 LOEWER, Barry C. 2007. "Mental Causation, or Something Near Enough." in *Con-*
2928 *temporary Debates in Philosophy of Mind*, edited by Brian P. MCLAUGHLIN and

- Jonathan COHEN, pp. 243–264. *Contemporary Debates in Philosophy* n. 8. Oxford: Blackwell Publishers. Second edition: McLaughlin and Cohen (2023).
- MARSHALL, Dan. 2021. “Causation and Fact Granularity.” *Synthese* 199(3-4): 8029–8045, doi:10.1007/s11229-021-03151-2.
- MCDERMOTT, Michael. 1995. “Redundant Causation.” *The British Journal for the Philosophy of Science* 46(4): 523–544, doi:10.1093/bjps/46.4.523.
- MCKAY, Thomas J. and NELSON, Michael. 2010. “Propositional Attitude Reports.” in *The Stanford Encyclopedia of Philosophy*. Stanford, California: The Metaphysics Research Lab, Center for the Study of Language and Information, <https://plato.stanford.edu/archives/win2010/entries/prop-attitude-reports/>.
- MCLAUGHLIN, Brian P. and COHEN, Jonathan, eds. 2007. *Contemporary Debates in Philosophy of Mind*. Contemporary Debates in Philosophy n. 8. Oxford: Blackwell Publishers. Second edition: McLaughlin and Cohen (2023).
- , eds. 2023. *Contemporary Debates in Philosophy of Mind*. 2nd ed. Contemporary Debates in Philosophy n. 8. Hoboken, New Jersey: Wiley-Blackwell. First edition: McLaughlin and Cohen (2007), doi:10.1002/9781394259847.
- MILLER, Alexander. 1995. “Objectivity Disfigured: Mark Johnston’s Missing-Explanation Argument.” *Philosophy and Phenomenological Research* 55(4): 857–868, doi:10.2307/2108336.
- MOLNAR, George. 2003. *Powers: A Study in Metaphysics*. Oxford: Oxford University Press. Edited by Stephen Mumford, doi:10.1093/acprof:oso/9780199204175.001.0001.
- MOORE, Dwayne. 2018. “Mind and the Causal Exclusion Problem.” in *Internet Encyclopedia of Philosophy*. University of Tennessee at Martin, <https://iep.utm.edu/mind-and-the-causal-exclusion-problem/>.
- PARTEE, Barbara Hall. 1989. “Possible Worlds in Model-Theoretic Semantics: A Linguistic Perspective.” in *Possible Worlds in Humanities, Arts and Sciences*, edited by Sture ALLÉN, pp. 93–123. *Research in Text Theory* n. 14. Berlin: de Gruyter. Proceedings of Nobel Symposium 65, doi:10.1515/9783110866858.93.
- PERRY, John R. 1989. “Possible Worlds and Subject Matter: Discussion of Partee (1989).” in *Possible Worlds in Humanities, Arts and Sciences*, edited by Sture ALLÉN, pp. 124–137. *Research in Text Theory* n. 14. Berlin: de Gruyter. Reprinted with a postscript in Perry (1993, 173–192) and in Perry (2000, 145–160), doi:10.1515/9783110866858.124.
- . 1993. *The Problem of the Essential Indexical and Other Essays*. New York: Oxford University Press, doi:10.1093/oso/9780195049992.001.0001.
- . 2000. *The Problem of the Essential Indexical and Other Essays*. Stanford, California: CSLI Publications. Enlarged edition of Perry (1993).
- SARTORIO, Carolina. 2005. “Causes as Difference-Makers.” *Philosophical Studies* 123(1/2): 71–96, doi:10.1007/s11098-004-5217-y.

- 2969 SOBER, Elliott R. 1982. "Why Logically Equivalent Predicates May Pick Out Different
2970 Properties." *American Philosophical Quarterly* 19(2): 183–189.
- 2971 STALNAKER, Robert C. 1984. *Inquiry*. Cambridge, Massachusetts: The MIT Press.
- 2972 WALSH, Sean Drysdale. 2016. "Predicativity, the Russell-Myhill Paradox, and Church's
2973 Intensional Logic." *Journal of Philosophical Logic* 45(3): 277–326, doi:[10.1007/s1](https://doi.org/10.1007/s10992-015-9375-5)
2974 [0992-015-9375-5](https://doi.org/10.1007/s10992-015-9375-5).
- 2975 YU, Andy Demfree. 2017. "A Modal Account of Propositions." *Dialectica* 71(4): 463–488,
2976 doi:[10.1111/1746-8361.12193](https://doi.org/10.1111/1746-8361.12193).

PROOF

PROOF

2977

Lewisian Worlds and Buridanian

Possibilia

2978

BOAZ FARADAY SCHUMAN

2979 Many things can be other than they are. Many other things cannot. We
2980 talk about such things all the time. But what is this talk about? One an-
2981 swer, presently dominant in analytical philosophy, is that we are speaking
2982 of possible worlds: if something can be other than it is, then it actually
2983 is that way in some (other) world. If something cannot be otherwise, it
2984 is not otherwise in any world whatsoever. But what are these worlds?
2985 David Lewis famously claims that every world exists, just like ours does.
2986 In contrast, the medieval thinker John Buridan understands modal logic
2987 in terms of objects and causal powers: if something can be other than
2988 it is, then there is a causal power that can make it that way. If it cannot,
2989 then no causal power—not even God—can make it otherwise. As we'll
2990 see, (i) the Lewisian plurality is not possible on Buridan's account, and
2991 accordingly (ii) a basic tenet of classical theism is untenable on Lewis's
2992 metaphysics. In short, either the Lewisian plurality is incoherent, or a
2993 core monotheistic tenet is impossible.

2994 Modal sentences deal with things that can or must or cannot be. For example,
2995 we say that a triangle *can* be drawn, *must* be three-sided, and *cannot* be round.
2996 What makes a modal sentence modal? Short answer: its inclusion of a modal
2997 term like *can* (*possibly*), *must* (*necessarily*), and so forth. Such terms register
2998 that a claim is being qualified in such a way that the conditions of its truth
2999 are not limited to the way things actually are. But what is this modal talk
3000 about? Over the past two and a half millennia, answers have varied. Relatively
3001 recently, we have come to think of modes in terms of quantification over
3002 worlds: what is possible is true in at least one world, and what is necessary is
3003 true in all. Call this the *worlds-reading* (WR) of modal sentences. David Lewis
3004 (1941–2001) famously understands WR ontologically: these worlds really exist
3005 as spatiotemporal isolates, and are every bit as real as our own.

3006 Contrast WR with a much older—and for a long time prominent—
 3007 understanding of what modes are: terms whose operation on sentences
 3008 expands (or *ampliates*) the extension of their terms, so that the terms range
 3009 over possible objects, including non-existent ones. The modal properties of
 3010 these objects are grounded in the causal powers of existing things: a triangle
 3011 can be drawn because you or I can draw one; it is necessarily three-sided
 3012 because there is no causal power (not even God) capable of making a triangle
 3013 to be otherwise—at least, not without depriving it of its triangularity. Call
 3014 this the *objects reading* (OR) of modal sentences. This is the view of John
 3015 Buridan (c.1300–1361).¹ A careful examination of WR and OR reveals that (i)
 3016 they are incompatible, so that the Lewisian plurality is not a possible object or
 3017 collection of objects; and accordingly that (ii) the worlds-reading, at least in
 3018 its Lewisian form, is incompatible with a basic tenet of classical monotheism.

3019 Why compare Buridan and Lewis? I have three reasons. First, Lewisian
 3020 modal realism is well-known, and therefore provides a convenient off-the-
 3021 shelf foil for Buridan's modal ontology. Second, Lewis has clear ontological
 3022 commitments, and so he is easy to pin down. Compare the ontologically
 3023 agnostic Kripkean modal semantics and syntax: you and I may have very
 3024 different views on what worlds are, but nevertheless agree on a Kripkean
 3025 reading of the claims of WR.² So the Kripkean account does not provide a
 3026 clear and illuminating contrast for Buridan's modal ontology, the way Lewis's
 3027 approach does.³ Third, contrasting Lewis and Buridan illuminates latent
 3028 aspects of both. It gives us an insight into Lewis, hitherto unrecognised in
 3029 the literature; and it reveals Buridan's own views on the limitations on divine
 3030 power—limitations he does not explicitly discuss at length. After all, placing
 3031 restrictions on God's power would have been a hazardous thing to do at the
 3032 fourteenth-century University of Paris.⁴ All the more so for an Arts Master

1 For a discussion of earlier debates about causal powers in the twelfth and thirteenth centuries, see Peter King (2021).

2 This is indeed one of the strengths of modern modal logics: we can use them to talk about worlds, computation states, etc.—as Patrick Blackburn, Johan van Benthem, and Frank Wolter point out in the preface to their (2007) *Handbook of Modal Logic*. More on this in a moment.

3 I do not doubt the value of contrasting Buridan's account with other views of what worlds are: Adams's (1974) reduction of modal talk to sets of propositions, Plantinga's (1974) view that they are states of affairs, Stalnaker's (1976) that they are properties, or Fine's (2003) that they are facts. I suspect that for various reasons, some subtle and some less so, Buridan's anti-realism would preclude all these. But Lewis's view is the most unlike Buridan's and therefore provides a useful foil for now.

4 In particular the infamous Condemnations of 1277 insisted on the boundlessness of divine power. For a discussion, see Grant (1979), and more recently Thijssen (2018).

3033 who, as he explicitly acknowledges, is not qualified to teach theology.⁵ All the
 3034 same, we can tease out the consequences of the views Buridan does express.
 3035 And there is more here than meets the eye.

3036 Let's begin with WR, which is relatively familiar, and has two important
 3037 shortcomings that point to two strengths of OR.

3038 **1 Possible Worlds**

3039 Nowadays, we tend to think of modality in quantificational terms: a modal is
 3040 a sentence with a modal operator like '□' or '◇', for necessity and possibility,
 3041 respectively. Such operators quantify across possible worlds. On these lights,
 3042 □ φ just says that φ holds in all possible worlds, and ◇ φ says that φ holds in
 3043 at least one. The parallel, then, is with the ordinary first-order quantifiers:
 3044 (□-like) '∀', and (◇-like) '∃'.⁶

3045 There is much to be said for WR, but here I will limit myself to two points.
 3046 First, it's versatile: we can use the apparatus of worlds to construct a wide
 3047 variety of systems of alethic modal logic—that is, modal systems dealing
 3048 with necessary truths, possible truths, and so on. We can characterise an
 3049 astonishing number of systems in this way, and haggle about which one is
 3050 best (or best for what). We can also characterise non-alethic systems to model
 3051 knowledge and belief (epistemic logic), past, present, and future time (tense
 3052 logic), and morality (deontic logic). WR, then, is extremely fruitful.⁷

3053 Second, the WR is precise: can we give clear quantificational definitions of
 3054 terms like *necessarily* and *possibly*, which might otherwise seem qualitative
 3055 and murky. And, using Kripke's apparatus of frames, we can characterise our

5 That Buridan never advanced beyond the post of arts master, and so—in spite of his evident brilliance—never taught at the higher and more prestigious Faculty of Theology, is remarkable. In modern terms, this would be a bit like deciding to remain an assistant professor for life, even when promotion was available. For a discussion, see Jack Zupko (2003, xi–xii).

6 One need not, however, be committed to a semantics of possible worlds in order to think of modal terms quantificationally: already in 1924, well before the possible-worlds innovations of Kripke, Otto Jespersen pointed out that “necessity means that *all* possibilities are comprised, just as impossibility means the exclusion of all possibilities” (1924, 325, emphasis original). And before him, Avicenna (980–1037) gave a temporal reading of the modal operators, where what is necessary is always true, and always true is linked with the quantificational term every. See Khaled El-Rouayheb (2009, 210).

7 As Graham Priest (2016, 2653) puts it, “the clarity of the mathematics involved, and their usefulness in an analysis of many things other than modality—such as conditionals, meaning, knowledge and belief—meant that they [i.e., possible worlds] soon became part of the intellectual landscape.”

3056 systems with mathematical precision. But beyond all this, we might wonder:
 3057 what are these worlds, anyway?

3058 *Lewisian Worlds*

3059 David Lewis's answer to this question is famous and bold: all possible worlds
 3060 exist, and they are just as real as ours. As he tells us (1986, 2):

3061 The other worlds are of a kind with this world of ours. To be
 3062 sure, there are differences of kind between things that are parts of
 3063 different worlds [...] but [...] the difference between this and the
 3064 other worlds is not a categorical difference. Nor does this world
 3065 differ from the others in its manner of existing.

3066 So there are many worlds—as many, in fact, as there are ways things can be.
 3067 This ontological account of WR prompts two questions: how are these worlds
 3068 externally distinct from each other, and how are they internally unified?
 3069 Answers to both questions turn on spatiotemporal relations. To the former,
 3070 Lewis tells us (1986, 3):

3071 There are no spatiotemporal relations at all between things that
 3072 belong to different worlds. Nor does anything that happens at one
 3073 world cause anything to happen at another. Nor do they overlap;
 3074 they have no parts in common.

3075 Lewis frequently treats causation as the paradigmatic spatiotemporal relation.
 3076 Since the worlds have no spatiotemporal relations to one another, *a fortiori*,
 3077 there can be no causal interactions between them. They are, therefore, em-
 3078 phatically not like planets that are too far removed to interact with each other.
 3079 Instead, they are spatiotemporal isolates. Call this Lewis's *isolation doctrine*.

3080 Importantly, Lewis does not say that different worlds *cannot* interact, as
 3081 if blocked from doing so. Rather, they just *do not*: the notion of interaction
 3082 between different worlds makes no sense within his theory. This requirement
 3083 has a stipulative flavour—and, indeed, it is precisely that: a stipulation. This
 3084 point is important, and we will return to it in section 3.

3085 In like manner, Lewis accounts for the unity of worlds in terms of spa-
 3086 tiotemporal relations (1986, 71):

3087 If two things are spatiotemporally related, they are worldmates
 3088 [...] things are worldmates iff they are spatiotemporally related.

3089 A world is unified, then, by the spatiotemporal interrelation of its
 3090 parts.

3091 Again, this is presented in a stipulative way, though it is a corollary of the
 3092 doctrine of isolation: worlds are spatiotemporally isolated, and therefore
 3093 any spatiotemporally related things belong, *eo ipso*, to the same world. Here,
 3094 whether or not causal interaction *actually* occurs is less important than imme-
 3095 diately above: there does not need to be any obvious causal relation between
 3096 two things for them to belong to the same world. A long-dead star too distant
 3097 from Earth to interact with it nevertheless has spatiotemporal relations to us:
 3098 it is some distance away in time and space, and it came into being at some
 3099 time relative to us. It is, therefore, our worldmate.

3100 The foregoing considerations can be distilled into a precise account of
 3101 Lewisian worlds or *possibilia*, to wit:

3102 POSSIBILIA_L. A world *w* is an isolated unity of spatiotemporally
 3103 interrelated parts. If *x* and *y* have any spatiotemporal relations, they
 3104 are members of the same world.

3105 Importantly, the spatiotemporal relation is Euclidean. Let *R* be the spatiotem-
 3106 poral relation, so that *Rxy* says that *x* is spatiotemporally (though not neces-
 3107 sarily causally) related to *y*. Then, by POSSIBILIA_L,

$$\forall xyz(Rxy \wedge Rxz \rightarrow Ryz)$$

3108 For clarity, we can also represent this diagrammatically, as follows:

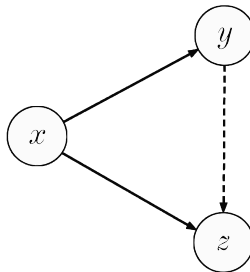


Figure 1: Euclidean R

3109 Here, R is represented by arrows; if the relation represented by the solid
 3110 arrows between x and y , x and z hold, then the relation represented by the
 3111 dotted arrow between y and z also holds.

3112 This fact makes the case that the Lewisian plurality is impossible (set out
 3113 in section 3) much easier to make, so let's linger on it for a moment. Let Rxy
 3114 and Rxz . It follows that Ryz . If it didn't, then x would be worldmates with
 3115 two objects that are not themselves worldmates with each other. So there
 3116 would be partial but incomplete overlap among at least two worlds. And
 3117 this goes against both $POSSIBILIA_L$, and against commonsense thinking about
 3118 spatiotemporal relations: if, for example, x is some spatial or temporal distance
 3119 from both y and z , then there must be some distance, however great, between
 3120 y and z themselves. Therefore, the spatiotemporal relation R is Euclidean.

3121 At the beginning of this section, I noted two significant advantages to
 3122 the WR of ordinary modal language: WR is precise, and fruitful. Before we
 3123 turn to the possible objects of Buridan, it's worth asking whether WR has
 3124 any drawbacks. For present purposes, I want to highlight two: WR does not
 3125 represent what is going on in ordinary modal language, and taken on its own,
 3126 it is uninformative about what grounds the modal properties of things.

3127 To begin with the latter: the extensional account furnished by WR does
 3128 not capture the ordinary notion of necessity *for* or *as*. For example, triangles
 3129 are necessarily three-sided; three-sidedness is necessary *for* triangle-hood.
 3130 Whereas you can paint a triangular object blue without removing its trian-
 3131 gularity, you cannot, say, rearrange its parts in such a way that it gains (or
 3132 loses) a side, and yet remains a triangle. This fact is not directly expressible
 3133 on WR; all it can tell us about this (or any other) necessary claim is that it
 3134 is true in every world. Fair enough, but such claims do not account for the
 3135 inseparability of three-sidedness and triangularity.

3136 Probably for this reason, most ordinary modal talk is not about worlds at all,
 3137 but rather about things, and the ways they can be in *this* world. Scott Soames
 3138 gives some remarks that support this point in his discussion of reference to
 3139 non-existent objects (2010, 128):

3140 Although this is controversial, the idea that we can refer to, and
 3141 quantify over, only things that exist is, I believe, an unfounded
 3142 philosophical prejudice at variance with our ordinary thought
 3143 and talk. For instance, imagine that I have all the materials to
 3144 build a doghouse, plus a plan specifying every detail of the design
 3145 and construction, including how each of the materials will be

3146 used. From studying the plan and materials, I know exactly which
 3147 structure I intend to create. Having identified it uniquely, I can
 3148 refer to it, predicate properties of it, and even name it.

3149 Soames's dog house is a possible, non-existent object. What makes it possible
 3150 is what *he* can do with materials and plans in *this* world. A lot of our day-to-
 3151 day modal talk is like this: when, for example, someone says they can paint
 3152 their house green, they are talking about *themselves*, and what they can do
 3153 with *their house*—not about their counterpart, in a relevantly similar world in
 3154 which their counterpart's house is green.

3155 Thus for all its versatility and precision, WR does not provide a full and
 3156 accurate report of what is going on in ordinary modal language. Such language,
 3157 judging by Soames's example, is about possible things, at least some of which
 3158 do not exist, whose modal properties are grounded in existing causal powers. I
 3159 have called this the objects reading (OR) of modal language; it is the approach
 3160 taken by John Buridan. It turns out that objects like Soames's doghouse are
 3161 precisely what Buridan has in mind in his analysis of *possibilia*.

3162 Possible Objects

3163 In the WR of modal language, modes operate on whole sentences, quantifying
 3164 over possible worlds. In contrast, Buridan's modal logic is not propositional
 3165 but *terminist*; he thinks of modes as acting on sentences' terms.⁸ Hence in his
 3166 treatment of modal semantics in *Tractatus de Consequentibus* (2.4), he tells us
 3167 that:

3168 A sentence (*propositio*) [...] about possibility has a subject term
 3169 that is amplified (*ampliatum*) by the modal term that follows
 3170 it, so that it stands (*ad supponendum*) not only for those things
 3171 which exist, but also for those things which *can* exist even though

8 While Buridan's *possibilia* have not received much attention, a good deal has been said already about Buridan's modal syntax and semantics. To date, the most thorough treatment of his syntax is chapter 9 of Paul Thom's (2003). And, following the concluding suggestions in G.E. Hughes's (1989), Catarina Dutilh Novaes (2007, 79–114) and Spencer Johnston (2015; 2017, 41–43) have given detailed analyses of Buridan's logic in terms of possible worlds. Gyula Klima, too, has remarked in his monumental translation of Buridan's *Summulae de Dialectica* that Buridan's modal semantics contains “effectively the gist of the idea of modern possible-worlds semantics” (2001b, 82, n.123).

they do not. Hence in this way it is true that air can come from water, although this is not true of any air that presently exists.⁹

Air from water is, as Paul Thom (2003, 170) has observed, a simple account of boiling. The water in this pot could boil; but since it is not boiling, it is not true of any actual air that it came from this water. Hence this water is possible—but not actual—air. Elsewhere, Buridan gives the appealing example of vinegar that could be produced from this wine, but will not, simply because you are going to drink the wine first (*de Caelo*, 1.25).¹⁰ These are the non-existent possible objects—or *possibilia*—to which the modal terms expand—or amplate—the terms of a sentence.¹¹

What are these non-existent *possibilia*?¹² Buridan deals with *possibilia* obliquely in his logic and metaphysics, and so we will have to reconstruct his view from these discussions. Here, I present three key passages: one dealing with necessity, one with impossibility, and the last with possibility. Approaching Buridan’s account of the *possibilia* from these three angles will allow us to build up a consistent and robust picture of his views on what they are.

9 “propositio [...] de possibili habet subiectum ampliatum per modum sequentem ipsum ad supponendum non solum pro his quae sunt sed etiam pro his quae possum esse quamvis non sint. Unde sic est verum quod aer potest fieri ex aqua, licet hoc non sit verum de aliquo aere qui est.” (Unless otherwise stated, all translations here are mine.) Note that Buridan is here talking about *divided* (roughly, *de re*) modals; he deals with *composite* (roughly, *de dicto*) modals elsewhere. Now, immediately below this passage, Buridan tells us that a modal sentence “B is possibly A” is equivalent to “What is or can be B can be A.” An anonymous reviewer for this journal has remarked on the connection with Williamson’s (2013, sec. 1.3) distinction between two readings of “possible stick”: the *predicative* reading (“x is a stick and x could have existed”), and the *attributive* reading (“x could have been a stick”). Buridan’s own account looks, *prima facie*, more like the predicative reading; but perhaps the two are not equivalent. At any rate, this question could form the basis of a stand-alone paper.

10 Cf. Aristotle’s cloak in *Peri Hermeneias* 9, which can be cut up, but may also simply wear out first (19a12–16).

11 For an overview of Buridan’s semantic doctrine of modal ampliation, and a case for it as one of his most significant contributions to the development of logic, see Zupko (2003, 67–70; 2018, sec. 4).

12 An anonymous reviewer for this journal has remarked that the common use of the term *possibilia* is for non-existent (possible) things, and does not extend to existing things as well. This is how I use it here, though it should be borne in mind that all *actualia* are, for Buridan, *possibilia* as well. After all, everything actual is possible.

2.3.1 Necessity in the Prior Analytics

3189 If S is necessarily P, then (by modal duality) it is not possible for S not to
 3190 be P—something Buridan explicitly acknowledges (see, e.g., *Tractatus de*
 3191 *Consequentibus* II, 5). Yet this analysis faces a problem. As Buridan asks in his
 3192 *Quaestiones super libros “Analyticorum Priorum”* (QAPr 1.25), what is the
 3193 modal status of the following sentence?

3194 (1) Humans are animals.

3195 Is (1) necessarily true? In *Prior Analytics* 1.9 (310a31), Aristotle clearly thinks
 3196 so. And indeed, (1) serves as a stock example of a necessary truth in medieval
 3197 logic.¹³ Yet (1) is falsifiable, since God could annihilate all human beings. As
 3198 Buridan tells us (QAPr 1.25, arg. 3):

3199 If it were supposed that [(1)] were not necessary, it would be
 3200 because God is capable of annihilating every human being. And
 3201 in such a case, no human would exist, and so no human would
 3202 be an animal.¹⁴

3203 For Buridan, all affirmative sentences, including universals, have existential
 3204 import, in contrast with negative sentences (both universal and particular),
 3205 which do not. Thus Buridan would reject the reading of (1) given by classical
 3206 FOL ($\forall x[\text{Human}(x) \rightarrow \text{Animal}(x)]$), which is capable of vacuous truth.¹⁵
 3207 Since there is no vacuous truth for affirmatives, (1) can be rendered false
 3208 by the annihilation of its subject matter. Therefore, since (1) is falsifiable, it
 3209 expresses a contingent truth.

3210 Nor is this sort of contingency limited to sentences which, like (1), are taken
 3211 from the natural sciences. It is a problem even for geometry:

3212 If this were so, then no claim of geometry would be necessary
 3213 either, since God can just as well annihilate all magnitudes as all
 3214 human beings. And then it would follow that geometry would

13 Along with “God exists” and “No human is a donkey.” Modern logical textbooks prefer mathematically-flavoured examples like “The set of primes is denumerable” and “ $a = a$.” The conventionalised role of these stock examples is clear.

14 “Item, si poneretur quod non esset necessaria, hoc esset pro tanto quia deus posset annihilare omnem hominem; ideo nullus homo esset, et sic nullus homo esset animal.”

15 I have discussed this aspect of Buridan’s logic, in connection with the traditional Square of Opposition, in Schuman (2022, 205–208).

3215 not be a science, which everyone would regard as false and unac-
 3216 ceptable. (*QAPr* 1.25, arg. 3)¹⁶

3217 God can annihilate everything with magnitude, and therefore magnitude itself.
 3218 If God were to do that, then all the affirmative claims of geometry would be
 3219 false, since the things they deal with would not exist. This is a consequence
 3220 of Buridan’s anti-realism, which extends even to the objects of mathematics
 3221 and geometry: if it so happened that there were no triangular arrangements
 3222 of matter, then there would be no triangles (though it would still be possible
 3223 to think and talk about them, like the roses of yesteryear). The same holds for
 3224 all other geometric and mathematical objects.

3225 Worse, even if God never gets that destructive, a crisis remains: the mere
 3226 fact that geometric claims *could* be falsified by an act of divine will entails that
 3227 these claims are contingent. If the truth of any claim is contingent, so is its
 3228 subject matter. Since the subject matter of any Aristotelian science (*scientia*)
 3229 must be necessary, it follows that even geometry is not a science. We can
 3230 expect the other sciences—with the obvious exception of theology—to fare
 3231 no better, given that God could annihilate their subject matter, too. So can
 3232 there be any science (apart from theology) at all?

3233 Buridan’s answer is yes: the claims of geometry (and of the other sciences)
 3234 are necessary, but their necessity is attenuated: they are not necessarily true
 3235 *simpliciter*. Rather, they are true “so long as” or “just when” (*de quando*)
 3236 the things their subject and predicate terms stand for exist.¹⁷ Assuming no
 3237 annihilation of their subject matter occurs, they will remain true—indeed,
 3238 *necessarily* true:

3239 Necessity “just when” (*de quando*) comes about from the fact that,
 3240 whenever the subject and predicate terms do stand for anything,
 3241 they stand for the same thing (I am here speaking of affirmative
 3242 sentences). And in this way I say that the following are necessary:
 3243 “Humans are animals,” or also “Horses are animals.” Indeed, even
 3244 “A rose is a flower” is necessary in this way, even if there are no
 3245 roses now. And although there is not a lunar eclipse happening

16 “si hoc obstaret, nulla propositio geometrica esset necessaria, cum deus ita possit annihilare omnes magnitudines, sicut omnes homines. Et tunc ultra sequeretur quod geometria non esset scientia, quod reputatur ab omnibus falsum et inconveniens.”

17 For an analysis of this kind of necessity, in particular in connection with tense and the necessity of the past, see Normore (2013).

3246 right now, still the following is necessary: “A lunar eclipse is a
3247 lack of light [coming] from the sun.” (QAPr 1.25, co)¹⁸

3248 So a sentence like (1) is necessarily true, assuming the existence of the things
3249 it deals with, namely, humans. Likewise, the claims of astronomy are true
3250 even when the events they describe are not presently occurring, since any
3251 time they *do* occur, the sentences are true. Thus, according to the account
3252 set out by Buridan in QAPr 1.25, a sentence like (1) can only be falsified by
3253 the *annihilation* of the things it deals with. There is no way to falsify (1) that
3254 leaves humans intact. So whenever humans exist, (1) is true.

3255 Thus the contrast between necessity and contingency in terms of modality
3256 simply construed (*simpliciter*) is the contrast between unfalsifiability and
3257 falsifiability. The contrast between necessity and contingency in terms of *de*
3258 *quando* modality is the contrast between falsifiability only by annihilation
3259 (*de quando* necessity) and falsifiability by alteration (*de quando* contingency).
3260 That humans are animals is *de quando* necessary, because it can only be rendered
3261 false by the removal of its subject matter—that is, by annihilating all
3262 humans. In contrast, the fact that some humans are bearded is *de quando* con-
3263 tingent, since shaving them alters the fact, but leaves the subjects essentially
3264 intact.

3265 From these observations, we can give the following Buridanian definition
3266 of necessity:

3267 BURIDANIAN NECESSITY. S is necessarily P just in case S can only
3268 be made to be not-P by annihilating S.

3269 This provides a good starting point for an account of Buridanian modality;
3270 however, there are crucial ambiguities that must be sorted out if the above
3271 definition is to be consistent with the others we will look at below. Its adoption
3272 here is, therefore, tentative.

18 “Necessitas de quando ex hoc provenit quod oportet subiectum et praedicatum quodcumque supponunt pro aliquo supponere pro eodem; et hoc dico in affirmativis. Et sic dico quod haec est necessaria ‘homo est animal,’ vel etiam ‘equus est animal.’ Immo etiam haec est necessaria ‘rosa est flos,’ licet modo nulla sit rosa. Et quamvis non sit eclipsis lunae, tamen haec est necessaria ‘eclipsis lunae est defectus luminis a sole.’”

2.2 *Impossibility in the Peri Hermeneias*

3274 In *Peri Hermeneias* 2 (16a19), Aristotle tells us that nouns (ὀνόματα; Aristoteles
3275 Latinus: *nomina*) have signification. But Buridan asks, what about nouns like
3276 *chimera*, which do not signify anything at all?

3277 It is asked: does every noun (*nomen*) signify something?

3278 Objection: it does not, because the term *chimera* signifies nothing
3279 apart from a chimera. And yet a chimera is nothing. Therefore, it
3280 signifies nothing whatsoever. (*Peri. Herm.* 1.2, arg. 1)¹⁹

3281 A chimera not only does not exist, like the roses of yesteryear; it is, in fact,
3282 impossible. Buridan makes this point several times: the chimera is made of
3283 impossible parts.²⁰ In this respect, we may take it to be just like Schopen-
3284 hauer’s wooden iron or Frege’s square circle (Schopenhauer 1819, vols. 1,
3285 para.53; Frege 1884, para. 74). Because the chimera cannot exist, it cannot be
3286 signified. And this seems to present a semantic counterexample to the *Peri*
3287 *Hermeneias* definition of nouns, even though syntactically, *chimera* functions
3288 like any other noun.

3289 Buridan’s solution here is to treat *chimera* as equivalent with the phrase
3290 “animal made up of parts that cannot be combined,” and to note that, although
3291 this whole phrase does not signify anything, it has significative parts (namely,
3292 *animal* and *part*). The details of this solution need not detain us here. What is
3293 significant for our purposes is the role of the chimera as an impossible object,
3294 whose impossibility is a function of its putative combination of impossible
3295 parts. We can use such *impossibilia* for our next definition:

3296 BURIDANIAN IMPOSSIBILITY. S is not possibly P if S and P cannot
3297 be combined.

3298 This relatively straightforward definition will be useful for untangling an
3299 important exegetical problem in section 2.4.

19 “queritur utrum omne nomen significat aliquid. Arguitur quod non, quia iste terminus ‘chimaera’ nihil significat aliud a chimaera. Et tamen nihil est chimaera. Ergo nihil omnino significat.”

20 “chimaera est animal compositum ex membris ex quibus impossibile est aliquod animal componi” (*De Demonstrationibus* 8.2.3). For a lively discussion of the role of the chimaera in the history of philosophy, see Ebbesen (1986).

203 Potency in the Metaphysics

3301 In his *Questions on the “Metaphysics” of Aristotle (QM)* 9.5, Buridan raises the
 3302 question of whether everything that something *will* do can be said to be what
 3303 it is *able* to do. If so, we get some strange results, as Buridan points out:

3304 A horse can come from wool. For earth comes from wool [by
 3305 decomposition], and grass come from the earth, and from those
 3306 grass which perhaps a horse will eat there can come horse semen,
 3307 and, at length, another horse. And so even a horse can come from
 3308 wool. And the same holds for all other modes of transmutation.²¹

3309 Here the problem is whether or not the relation between S and P expressed
 3310 by “S is possibly P” is transitive: if S can be P, and P can be Q, does it follow
 3311 that S can be Q?

3312 No, says Buridan: when we say that S can be P, we are generally speaking
 3313 in terms of a *proximate* potency, rather than a remote one: S is proximately
 3314 possibly P if S can become P in no more than one transmutation. In this way,
 3315 wool is possibly earth, because it can become earth in one transmutation (i.e.,
 3316 decay); similarly, earth can become grass, and so on. Any other potencies that
 3317 require multiple transmutations are remote—as is, for instance, the potency
 3318 of wool to become a horse. Hence Buridan tells us that:

3319 Aristotle concludes the opposite. For he asks, when should some-
 3320 thing be said to be in potency, and when should it not? And he says
 3321 that something should not be said to be in potency with respect
 3322 to some form, except when only one transmutation is required,
 3323 by which that form may be imparted on it. (*QM* 9.5, fol. 58rb)²²

3324 So although remote potencies can be discussed transitively, proximate poten-
 3325 cies cannot. If the two are conflated, as in the horse-from-wool example, then,

21 “Similiter ex eadem lana potest fieri equus, quia ex lana fiet terra, de inde herba, et ex illa herba forte quam equus comedet poterit fieri sperma equi et tandem equus. Et ita etiam ex lana potest fieri equus. Et sic de omnibus aliis modis transmutandi” (*QM* 9.5, fol. 58rb). Among the other modes of transmutation Buridan discusses here are “Wool can become a hatchet” (wool > earth > stone > iron > hatchet), and “An infant can build a house” (infant > adult human > carpenter).

22 “Oppositum determinat Aristoteles. Querit enim quando aliquid debeat dici in potentia et quando non. Et dicit quod aliquid non debet dici in potentia ad aliquam formam, nisi quando sola transmutatio requiritur per quam illa forma perducatur.” Buridan seems to have in mind Aristotle’s *Physics* 1.4 (188a32–b3).

3326 according to Buridan, the result is an equivocation.²³ Thus, although wool
 3327 can decompose into earth, grass can grow from earth, and so forth, it does not
 3328 follow that wool can become grass—much less a horse. Hence in speaking of
 3329 possible horses, we are not speaking of all the things that, through multiple
 3330 transmutations, could become a horse. If we were, then everything would
 3331 be a possible horse, since, as Buridan observes, “anything can come from
 3332 anything—albeit through several transmutations” (*QM* 9.5, fol. 58rb).²⁴

3333 So much for *possibilia* arising from natural causes, like possible dirt that
 3334 can be generated from wool. But a problem remains: why couldn’t God just
 3335 rearrange the matter in a horse, say, to make it into a pile of dirt? So then a
 3336 horse is possibly dirt (and vice-versa).²⁵ And if so, then our main problem
 3337 comes roaring back: everything is possibly everything.

3338 Buridan himself does not consider this problem, but there is indirect textual
 3339 evidence that he would reject such a claim: after all, he frequently tells us that
 3340 the following is impossible:

3341 (2) A human is a donkey.

3342 Granted, it is not beyond divine power to transform the matter of a human
 3343 being into a donkey by imparting on it the appropriate form. But again, (2) is
 3344 impossible. How?

3345 The solution is to appeal to the notion of change entailing annihilation (or
 3346 destruction—more on this in a moment), which we saw above in connection
 3347 with *de quando* necessity. For example, consider the following sentence:

3348 (3) Socrates is a human.

3349 Any formulation of (3) is true whenever Socrates exists. And while (3) can be
 3350 rendered false, this can only happen by the destruction of Socrates. Similarly
 3351 if, instead of being served a hemlock cocktail, Socrates met his demise by
 3352 having his matter suddenly morphed into the form of a donkey, (3) would

23 “Modo in proposito est bene aequivocatio de potentia propinqua et remota” (*QM* 9.5, fol. 58va).

24 “quia ex quolibet potest fieri quodlibet—licet per multas transmutationes.”

25 Given the essentialist background of the current discussion, I’m aware I am treading dangerously close to an old problem at which even young Socrates is reported to have balked: Does dirt have an essence? (*Parmenides* 130c–d). I wish to remain neutral on this point: for my purposes, the only concession I have to make is that whatever makes horses horsey is essentially different from whatever makes dirt dirty. Maybe I beg the question on this. But I invite you to beg it with me. After all, we’re in good company, historically speaking.

3353 become false. But so would the claim that Socrates himself is a donkey, since
 3354 Socrates himself would no longer exist. So Socrates is not possibly a donkey.

3355 We have limited ourselves to transmutation in talking about things-possibly-
 3356 being-other-things, and to one transmutation at that. Granted, then, God can
 3357 morph Socrates's matter into a donkey. But this morphing does not count as a
 3358 transmutation in the natural sense, nor is it a potency belonging to Socrates.
 3359 And so this fact no more entails that Socrates is a possible donkey than does
 3360 the fact that Socrates can die and decay into soil, which then nourishes a
 3361 plant, which a donkey eats, etc.

3362 Here, then, we return to the original claim that *impossibilia* are impossi-
 3363 ble combinations: donkey-Socrates, chimaeras—anything, in short, made
 3364 up of parts that cannot be combined. Soon, we will see that Lewisian possible
 3365 worlds, too, are Buridanian *impossibilia*. But first, we have to find a way of
 3366 making the foregoing definitions consistent.

264 What Are Buridanian Possibilia?

3368 In a seminal (1989) paper, G.E. Hughes raises several questions about Buri-
 3369 dan's modal logic and its underlying ontology. Concerning the latter, he tells
 3370 us (1989, 97):

3371 For a long time I was puzzled about what Buridan could mean by
 3372 talking about possible but non-actual things of a certain kind. Did
 3373 he mean by a "possible A," I wondered, an actual object which is
 3374 not in fact A but might have been, or might become, A? My house,
 3375 e.g., is in this sense a possible green thing because, although it
 3376 is not in fact green, it could become green by being painted. But
 3377 this interpretation won't do; for Buridan wants to talk, e.g., about
 3378 possible horses; and it seems quite clear that he does not believe
 3379 that there are, or even could be, things which are not in fact horses
 3380 but which might become horses.

3381 Notably, Hughes makes no mention of the *Metaphysics* discussion—about
 3382 horses, too!—which we just considered. This comes as no great surprise: that
 3383 text is, to this day, neither edited nor translated.²⁶

26 Granted, Hughes himself did know Latin and was experienced in palaeography. He even edited a portion of the *Logica Magna* of Paul of Venice (ca. 1369–1429). Still, one can't read everything.

3384 Even so, Hughes's initial proposal is actually quite close to Buridan's own
 3385 account: a house is a possible green thing, because there are powers in the
 3386 world capable of making it so. The issue of substantial change—things be-
 3387 coming horses—is somewhat more thorny, since it seems a little odd to speak
 3388 of things which are not horses, but which could become horses, as Hughes
 3389 observes. And yet this is precisely what we are warranted to do, as Buridan
 3390 explicitly tells us, provided we limit ourselves to at most one transmutation:
 3391 horse semen is not a horse, but it is a possible horse.

3392 Frustrated by his version of the horse puzzle, and unaware of Buridan's *QM*
 3393 discussion, Hughes falls back on the familiar framework of possible worlds
 3394 to describe what Buridan is doing:

3395 What I want to suggest here, very briefly, is that we might under-
 3396 stand what he says in terms of modern "possible world semantics."
 3397 Possible world theorists are quite accustomed to talking about
 3398 possible worlds in which there are more horses than there are in
 3399 the actual world. And then, if Buridan assures us that by "Every
 3400 horse can sleep" he means "Everything that is or can be a horse
 3401 can sleep," we could understand this to mean that for everything
 3402 that is a horse in any possible world, there is a (perhaps other)
 3403 possible world in which it is asleep. It seems to me, in fact, that in
 3404 his modal logic he is implicitly working with a kind of possible
 3405 worlds semantics throughout. (Hughes 1989, 97)

3406 Here, Hughes first claims that Buridan's modal logic can be understood using
 3407 the modern apparatus of possible worlds semantics. But then he strengthens
 3408 that claim: Buridan is in fact working with possible-worlds semantics, however
 3409 implicitly.

3410 From what we've seen of Buridan so far, we can tell that at least the latter
 3411 claim is mistaken. Buridan's view of modality is grounded in *causation*: if
 3412 there exists no power to make S to be not P (at least without annihilating S),
 3413 then S is necessarily P. Likewise, if S can be made to be P (through at most one
 3414 transmutation), then S is possibly P. Thus something's modal properties are
 3415 grounded in the powers that exist *in this world*, which are capable of making
 3416 it to be this or that way. In other words, Buridanian *possibilia* are, in general
 3417 terms, objects, some of them nonexistent, whose modality depends on the
 3418 causal powers of actually existing things. Since one of these existing things is
 3419 the Almighty, and since the Almighty exists by simple (which is to say strictly

3420 unalterable) necessity, the modal properties of the *possibilia* are stable. There
3421 are no other worlds in the picture.

3422 So much for what Buridan's view is not. But the definitions we've distilled
3423 from the texts face an important exegetical problem: both necessity, on one
3424 hand, and possibility, on the other, are each in their own way inconsistent with
3425 the account of impossibility as sketched above. Impossibility, unlike necessity,
3426 does not turn on annihilation: a chimaera is made up of impossible parts,
3427 not parts that would be literally reduced to nothing if they were combined.
3428 Moreover, there are diachronic possibilities, such as a human turning into a
3429 corpse, which are not synchronically possible: a human cannot be inanimate
3430 and rational at the same time, even though every human is a possible corpse.
3431 Just like *chimera*, *inanimate rational animal* therefore picks out an impossible
3432 object. The language of transmutations is therefore not applicable to syn-
3433 chronic impossibilities. These facts call for a re-examination of necessity
3434 and of possibility as set out above. We will soon see that (i) these accounts can,
3435 happily, be made consistent, and (ii) that the consistent account that emerges
3436 gives us a straightforward definition of Buridanian *possibilia*.

3437 First, the account of necessity, which turns on annihilation (rather than
3438 destruction) of the subject is too strong. There is more than one way to
3439 make Socrates not a human, and thereby to deprive him of his essence: one
3440 is through (divine) annihilation—literal reduction to nothing. Another is
3441 through (divine or natural) destruction—undergoing a change that entails
3442 removal of his (human) essence. After all, following his death, Socrates is no
3443 longer a human, but this fact does not turn on any annihilation of Socrates.

3444 Why then does Buridan discuss necessity in terms of annihilation at all? Re-
3445 call that, in the *QAPr*, Buridan is (*inter alia*) worried about the falsification of
3446 geometrical claims: if *all* magnitudes were annihilated, then the propositions
3447 of geometry would be rendered false. But this would not follow if everything
3448 with mass were simply destroyed—that is, if everything now existing were
3449 reduced to an undifferentiated soup. Even in that soup, there would be at least
3450 some dimension, surface, and so on to speak of. Conversely, the claim that
3451 humans are animals *would* be falsified if all humans were destroyed—that is,
3452 if everyone died all at once. Hence it seems that the reliance on annihilation
3453 is stronger than it needs to be for the definition of humans as animals, though
3454 perhaps not for the propositions of geometry taken collectively. I therefore
3455 propose a weakening of this requirement, at least for our definition of *possi-*
3456 *bilia*: S is necessarily P, just in case S cannot be made other than P without
3457 *destroying S*.

3458 The second exegetical problem is that the definition of possibility is quite
 3459 weak: supposing that S is possibly P just in case S can become P through at
 3460 most one transmutation, it follows that Socrates, while still alive, is possibly a
 3461 corpse. Fair enough; but, as we observed, the combination of Socrates, *qua*
 3462 rational animal, and corpse, *qua* inanimate object, is impossible.²⁷ Therefore,
 3463 the most straightforward reading of impossibility, set out in section 2.2, clashes
 3464 with the weak sort of possibility set out in section 2.3. What do we do?

3465 It is true that Socrates is possibly a corpse. And it is also true that Socrates,
 3466 while alive and barbate, is possibly clean-shaven. In the former case, Socrates
 3467 loses his essence; in the latter he does not. We should therefore distinguish
 3468 two kinds of change: one which involves loss of essence, but only through
 3469 one transmutation; and another which leaves the subject intact.

3470 Which kind of possibility is relevant to our purposes? *Impossibilia* are
 3471 impossible combinations; *possibilia* then should be possible ones. Since at
 3472 least some transmutations involve change into something impossible with
 3473 the essence of the subject, as our example of *rational animal* and *inanimate*
 3474 *object* shows, *possibilia* cannot include contrary diachronic states considered
 3475 synchronically. We should, therefore, take the stronger reading of possibility,
 3476 suggested by the account of impossibility: S is possibly P iff S can be P in a
 3477 way that does not entail the destruction of S.

3478 From these considerations, we can give the following definition of *possibilia*,
 3479 which balances out the accounts in Buridan's texts:

3480 POSSIBILIA_B. S is possibly P just in case there is a power to make S
 3481 to be P without destroying the essence of S.²⁸

3482 This definition casts a pretty wide net: *possibilia* will include not just the
 3483 various natural kinds and subkinds we see in the world, but also anything else
 3484 which could be produced by any power—including God—without destruc-

27 For a discussion of related problems in the logic and semantics of the twelfth century, see Cameron, M.A. (2015).

28 As an anonymous reviewer for this journal has pointed out, this definition, and the intuitions that motivate it, rest on essentialist assumptions. That is true, but the assumptions are weak ones: we need not assume that we have correctly identified the essence of S; we need only say that as a member of a natural kind, S has an essence—whether or not we know what it is. Still, one might worry about possibilities for houses and other artifacts, since (at least in Aristotelian metaphysics) artifacts do not have essences. A house, then, is possibly green, and also possibly a heap of rubble, and neither of these changes involves a loss of essence. Perhaps we could appeal to the house's function, which is preserved in the case of painting, but lost when it is reduced to rubble. But I leave that for another day.

3485 tion of the subject. So horses larger than planets are, presumably, (divinely)
 3486 possible; as are humans capable of walking on water, virgin mothers, and so
 3487 on. But conspicuously absent from this jungle of *possibilia* is the Lewisian
 3488 plurality of worlds with which we began.

3485 **3 Are Lewisian Possible Worlds Possible?**

3490 —Or, to put the question in Buridanian terms: can God create a Lewisian
 3491 plurality of worlds? First, the argument pro: it seems that God can indeed
 3492 create as many worlds as God pleases. Recall our account of the unity of
 3493 Lewisian worlds, set out above (section 1.1). So long as we conceive of a world
 3494 as just a cluster of spatiotemporally interrelated *possibilia*, there seems to be
 3495 no barrier in principle to clustering them. Here is why: some—and probably
 3496 most—possible objects are made up of interrelated possible parts. Consider,
 3497 for example, a possible watch that does not now exist. Such a possible watch
 3498 will not be undifferentiated all the way through, like liver paste, but will have
 3499 interrelated possible parts—possible gears, possible springs, etc.

3500 Now it would be arbitrary and just plain wrong to place a limit on how
 3501 large such a possible object could be, at least in terms of what God can create:
 3502 if a watch can be made the size of a tower clock, why not a watch the size of
 3503 Manhattan? Likewise, it would be arbitrary to place a limit on their complexity:
 3504 if a watch the size of Manhattan is permissible, why not a huge and complex
 3505 astronomical horologium—one as large and complex as our universe, even?

3506 From these considerations, we can distill two principles, namely:

- 3507 (i) *possibilia* can be internally complex, comprising interrelated possible
- 3508 parts; and
- 3509 (ii) there is no limit in principle to the size or complexity of such *possibilia*.

3510 From (i) and (ii)—so the argument runs—it follows that God could make
 3511 worlds, roughly construed as manifolds of interrelated objects.

3512 In fact, we can strengthen this claim: the *possibilia* just *have* to be in some
 3513 possible world. Consider a possible object, say a fork: Can such an object exist
 3514 outside a world or manifold? Or must any such possible object exist within
 3515 some kind of manifold? The existence of a fork outside some spatiotemporal
 3516 manifold seems, if not impossible, then at least a little weird. A fork in the
 3517 absence of other objects is one thing, but a fork in the absence of time and
 3518 space is quite another. And so, it seems, possible objects only ever inhabit

3519 possible worlds—a claim Jonathan Schaffer describes as a “platitude” (2010,
 3520 65). Thus a metaphysics of possible objects must, if it is to be coherent, collapse
 3521 into a metaphysics of possible worlds.²⁹ To put it in Buridanian terms: God
 3522 could create a fork, but not without creating it in some world.

3523 What this argument gets us is world-boundedness for anything God could
 3524 create. What it does not get us is a plurality of such worlds in a Buridanian
 3525 ontology. And there is good reason to reject such a plurality on these terms,
 3526 both for Buridan and for Lewis. Note first that these worlds are either actual,
 3527 in the sense that God has made them, or they are possible but non-existent, in
 3528 the sense that God has not made them, but could. In either case, the question
 3529 is: could God make an actual plurality of worlds? If so, then the Lewisian
 3530 plurality is possible; if not, then it is impossible.

3531 Following Lewisian doctrine, these worlds will have to be isolated: if they
 3532 are not, they no more count as distinct possible worlds than do planets in
 3533 different galaxies or, say, the city of Paris in different epochs. They must not
 3534 be at any spatiotemporal distance from each other. So can God create worlds
 3535 that are not worldmates in this way?

3536 Suppose God made these worlds. What does it mean to say such worlds
 3537 are causal isolates—i.e., that they cannot interact? Distance will not do the
 3538 trick: worlds are not causally isolated by any spatiotemporal distance, the way
 3539 you and I are isolated from a long-dead star in Andromeda. Space is not what
 3540 separates the worlds. Nor is time. Lewis has been clear.

3541 Perhaps we can say that God stipulates that the worlds cannot interact:
 3542 there is just an impermeable barrier between the worlds, analogous to the
 3543 glass plates separating different tanks in a divided aquarium, or the walls
 3544 splitting off different theaters in a cineplex. Perhaps it is physical, perhaps it
 3545 is by divine *fiat*. Either way, we then face three problems.

3546 First, what happens when two things in different worlds interact with the
 3547 dividing barrier or *fiat* that separates them? Suppose, for instance, that there
 3548 is a barrier between worlds Γ and Δ . Now γ and δ —possible objects in Γ and
 3549 Δ , respectively—are blocked from interacting by the barrier/*fiat* (imagine
 3550 fish bumping into the opposite sides of a glass aquarium divider). Then the
 3551 barrier that prohibits causal interaction between the two worlds, Γ and Δ ,
 3552 nevertheless causally interacts with both of them. Therefore, that barrier will
 3553 be a member of both worlds, according to Lewis’s definition: it has worldmates
 3554 on both sides. But preventing such world-straddling was precisely what the

29 I owe the gist of this argument to Douglas Campbell.

3555 barrier was supposed to do. Its very act of blocking makes it a world-straddler.
3556 We can try adding barriers so that the two barriers on the Γ and Δ sides are
3557 separated, a bit like parallel sheets of glass in a double-paned window. But
3558 then we get a regress: What keeps the barriers themselves apart? What would
3559 happen if one barrier collided with whatever separates it from the other? In
3560 any case, the barriers must both interact with whatever separates them.

3561 Second, even if God could somehow separate Γ and Δ causally from each
3562 other, it would still make sense to think of them as related temporally: just as
3563 we can speak of one movie in a cineplex starting at the midpoint of another,
3564 so we can speak of a universe being half as old as another—that is, as being
3565 created midway along the life cycle of another universe. For instance, we
3566 could reasonably ask whether, from God's perspective, the timeline of Δ is
3567 half as long as that of Γ , whether Δ already existed when Γ was created, and
3568 so on.

3569 Third, and most importantly, even if such worlds could be isolated from
3570 each other in a way that circumvents the foregoing two problems, they will still
3571 be causally related via their causal dependence on God. As we saw (section 1),
3572 the general spatiotemporal relation (though not necessarily causation) is
3573 Euclidean: if xRy and xRz , then zRy . Thus although two worlds may not
3574 causally interact, they are not spatiotemporally independent, since they have
3575 the same cause. They are, then, causal siblings, even if they never interact.
3576 And if they are produced by the same cause, then they are causally related, if
3577 only in virtue of being created by the same God.

3578 Lewis considers pseudo-pluralities (his term) like these, which, according
3579 to him, are not made up of truly isolated worlds. Their constituents are, rather,
3580 worldmates, even if locally they look like isolated worlds. Here is the one our
3581 cineplex and aquarium examples most closely resemble:

3582 The spacetime of the big world might have an extra dimension.
3583 The world-like parts might then be spread out along this extra
3584 dimension, like a stack of flatlands in three-space. (Lewis 1986,
3585 72)

3586 But, as Lewis is quick to point out, this is not a true plurality. For him, there is
3587 no way to speak of temporal relations across truly isolated worlds: if there is

anything like a common origin, then the worlds belong to the same manifold. And if they belong to the same manifold, they are not truly isolated.³⁰

Here is the most common objection I have faced to this line of reasoning: it is not that Lewisian worlds *cannot* interact, in the sense that there is some mechanism keeping them apart. Instead, they just *do not*. This is a restatement of Lewis's isolation doctrine. But the isolation doctrine is not a conclusion Lewis reaches by argument, as we saw in section 1. It is, rather, a stipulation. And in fact, this is how Lewis presents it: right up front, on the second page of his (1986) exposition. It is thus more a starting point than a destination.

Accordingly, no criticism of this doctrine can address Lewis's arguments for it, since he does not give us any. All that can be asked is whether it makes any sense. The answer, on Buridan's metaphysics (or any metaphysics that posits a single First Cause), is *No*. To anyone who espouses such a metaphysics, then, a Lewisian plurality of worlds must be something like Naive Set Theory: plausible on the face of it, but deep down self-contradictory. Lewis's worlds simply do not work on Buridan's framework. And, we might think, so much the better for Buridan.

I am not, by the way, the first to make any claims about the (in)compatibility of Lewisian worlds with classical theism, though the causal one I have been elaborating here is novel. Paul Sheehy (2006) sets out a number of problems for the classical theistic conception of God on Lewisian modal metaphysics. The most significant of these is his argument, suggested by Richard Davis (2008), that Lewisian possible worlds effectively chop God up, making each counterpart God a world-bound entity—an understanding that runs contrary to classical theism's commitment to divine unity. Ross Cameron, R.P. (2009) disagrees: Lewisian metaphysics can countenance abstracta existing outside of any world, as numbers do, so long as these abstracta are pure sets—that is, sets that contain only sets in their transitive closure (sets, sets-of-sets, sets-of-sets-of-sets, and so on, but no elements anywhere: only sets, including the empty set). God, it seems, could be such a set—even if it's doubtful whether such a set is what God's believers believe in (or, anyway, believe they believe in). Subsequent debate (Paul Sheehy 2009; Chad Vance 2016; Michael James Collier 2019, 2021) has dealt with this problem of divine (unitary) existence

30 Something similar could be said for the synchronic contrary possibilities of Scotus's (much discussed) *Lectura* I, dist. 39, q.1–5. Since these possibilities are rooted in the causal powers of a (single) will, they are worldmates. Therefore, these synchronic contrary possibilities are not true worlds in the Lewisian sense. For a discussion of Scotus in terms of possible worlds, see Wyatt (2000).

and world-boundedness, and whether, in these ways, God can be countenanced on Lewisian worlds. Brian Leftow (2012, 541–545) has, moreover, criticised Lewis on the grounds that positing one God is more economical than positing several (more on this in a moment).


For my part, I agree with Cameron and Collier that a Lewisian ontology can indeed countenance an abstract, un-world-bound Necessary Being of sorts. And I agree with Sheehy and Vance that Lewisian worlds are incompatible with classical monotheism, albeit for reasons apart from the ones they examine. After all, it is integral to classical theism that God has a creative—which is to say causal—role to play as well: God “created the heavens and the earth” (Genesis 1:1), is the One without Whom “nothing was made that was made” (John 1:2), the Originator, “Who commands only ‘Be!’ and it is” (*Al Baqarah* “The Heifer,” 117), and so on. (Countless other sources could be cited to this effect, but you get the idea). This central aspect of God’s activity is incompatible with Lewis’s doctrines about the plurality of worlds. Accordingly, possible worlds of the sort we have considered here will be deeply incompatible with (monotheistic) medieval philosophy in general—even if certain aspects of a given thinker’s modal logic or ontology might remind us of this (by now quite familiar) framework.³¹

What about Lewisian metaphysics considered in its own right? Even though a unified First Cause is not available on this framework, it does not follow that Lewis and his followers have to be atheists; if there is plurality in the worlds, there can also be a plurality of first causes. There is textual evidence that Lewis recognises this implication of his theory: in the introduction to the first volume of his (1983) *Philosophical Papers*, he remarks in passing that his view is consistent with the claim that “there are countless gods but none of them are our worldmates” (1983, xi). Since the worlds are, ontologically speaking, just like ours, it follows that our worldmates could include a local deity, and Lewis could merely be mistaken about the constituents of our actual world. So the Lewisian can still opt for a kind of polytheism, or mono-poly-theism, to

³¹ This will be true even when philosophical discussion centers on the notion of multiple worlds, e.g., in the claim of Al Ghazali and the Ash’arite theologians that God could have made other worlds than this one. Here, too, the worlds that could exist are referred back to a single unified power to bring them into existence, and so there is a similar problem for Lewis’s separation doctrine to the one discussed above. For a lively and interesting overview of this aspect of Al Ghazali’s thought, see Taneli Kukkonen (2000). (I am grateful to Silvia Di Vincenzo for bringing this to my attention.)

3651 adapt a term coined by Hart (2013, 127). But even basic classical monotheism
 3652 is, on these lights, impossible. For Lewisian ontology is a jealous god.*

3653 Boaz Faraday Schuman

3654  0000-0001-5763-8628

3655 University of Copenhagen

3656 boaz.schuman@hum.ku.dk

3657 References

- 3658 ADAMS, Robert Merrihew. 1974. "Theories of Actuality." *Noûs* 8(3): 211–231. Reprinted
 3659 in Loux (1979, 190–209), doi:10.2307/2214751.
- 3660 BLACKBURN, Patrick, VAN BENTHEM, Johan and WOLTER, Frank. 2007. *Handbook of*
 3661 *Modal Logic*. Studies in Logic and Practical Reasoning n. 3. Amsterdam: Elsevier
 3662 Science Publishers B.V.
- 3663 CAMERON, Margaret Anne. 2015. "The Logic of Dead Humans: Abelard and the
 3664 Transformation of the Porphyrian Tree." in *Oxford Studies in Medieval Philosophy*,
 3665 volume III, pp. 32–63. Oxford: Oxford University Press, doi:10.1093/acprof:oso/97
 3666 80198743798.003.0002.
- 3667 CAMERON, Ross P. 2009. "God Exists at Every (Modal Realist) World: Response to
 3668 Sheehy (2006)." *Religious Studies* 45(1): 95–100, doi:10.1017/S0034412508009827.
- 3669 COLLIER, Matthew James. 2019. "God's Necessity on Anselmian Theistic Genuine
 3670 Modal Realism." *Sophia. International Journal of Philosophy and Traditions* 58(3):
 3671 331–348, doi:10.1007/s11841-018-0659-4.
- 3672 —. 2021. "God's Place in the World." *International Journal for Philosophy of Religion*
 3673 89(1): 43–65, doi:10.1007/s11153-020-09764-w.
- 3674 DAVIS, Richard Brian. 2008. "God and Modal Concretism." *Philosophia Christi* 10(1):
 3675 57–74, doi:10.5840/pc20081014.
- 3676 DUTILH NOVAES, Catarina. 2007. *Formalizing Medieval Logical Theories: suppositio,*
 3677 *consequentiae and obligationes*. Logic, Epistemology, and the Unity of Science n. 7.
 3678 Dordrecht: Springer, doi:10.1007/978-1-4020-5853-0.
- 3679 EBBESEN, Sten. 1986. "The Chimera's Diary." in *The Logic of Being: Historical Studies*,
 3680 edited by Simo KNUUTTILA and Jaakko HINTIKKA, pp. 115–143. Synthese Histori-

* I am grateful for the insightful comments and criticisms of Irene Binini, Jon Bornholdt, Douglas Campbell, Nate Charlow, Jack Copeland, Graziana Ciola, Silvia Di Vincenzo, Peter King, Gyula Klima, Fred Kroon, Chris Martin, Calvin Normore, Diane Proudfoot, Jeremy Seligman, Trevor Teitel, and Andrew Withy. Thanks are also due to audiences at the University of Auckland, and the University of Canterbury, Christchurch, as well as audiences at meetings of the American Philosophical Association–Eastern, and the European Symposium of Medieval Logic and Semantics. The paper has, finally, benefitted from the feedback of two anonymous reviewers for *Dialectica*. Thanks, everyone!

- cal Library n. 28. Dordrecht: Kluwer Academic Publishers, doi:10.1007/978-94-009-4780-1_5.
- EL-ROUAYHEB, Khaled. 2009. "Impossible Antecedents and Their Consequences: Some Thirteenth-Century Arabic Discussions." *History and Philosophy of Logic* 30(3): 209–225, doi:10.1080/01445340802447905.
- FINE, Kit. 2003. "The Problem of Possibilia." in *The Oxford Handbook of Metaphysics*, edited by Michael J. LOUX and Dean W. ZIMMERMAN, pp. 161–179. Oxford Handbooks. Oxford: Oxford University Press. Reprinted in Fine (2005, 214–232), doi:10.1093/oxfordhb/9780199284221.001.0001.
- . 2005. *Modality and Tense: Philosophical Papers*. Oxford: Oxford University Press, doi:10.1093/0199278709.001.0001.
- FREGE, Gottlob. 1884. *Die Grundlagen der Arithmetik: Eine logisch-mathematische Untersuchung über den Begriff der Zahl*. Breslau: Wilhelm Koebner. Reissued as Frege (1961).
- . 1961. *Die Grundlagen der Arithmetik: Eine logisch-mathematische Untersuchung über den Begriff der Zahl*. Hildesheim: Georg Olms.
- GRANT, Edward. 1979. "The Condemnation of 1277, God's Absolute Power, and Physical Thought in the Late Middle Ages." *Viator* 10: 211–244, doi:10.1484/J.VIATOR.2.301526.
- HART, David Bentley. 2013. *The Experience of God: Being, Consciousness, Bliss*. New Haven, Connecticut: Yale University Press, doi:10.12987/9780300167337.
- HUGHES, George E. 1989. "The Modal Logic of John Buridan." in *Atti del convegno internazionale di storia della logica: le teorie delle modalità: San Gimignano, 5-8 dicembre 1987*, edited by Giovanna CORSI, Corrado MANGIONE, and Massimo MUGNAI, pp. 93–111. Bologna: Cooperativa Libraria Universitaria Editrice (CLUEB), <https://www.ontology.co/essays/hughes-buridan.pdf>.
- JESPERSEN, Otto. 1924. *The Philosophy of Grammar*. London: George Allen & Unwin, doi:10.4324/9780203716045.
- JOHN BURIDAN. 1964. *Quaestiones in Aristoteles Metaphysicam. Kommentar zur aristotelischen Metaphysik*. Frankfurt a.M.: Minerva.
- . 1976. *Tractatus de Consequentibus*. Philosophes médiévaux n. 16. Louvain la Neuve: Publications universitaires de Louvain. Edited by Hubert Hubien.
- . 1983. *Quaestiones longe super librum Perihermeneias*. Nijmegen: Ingenium Publishers. Edited with an introduction by Ria van der Lecq.
- . 1985. *Jean Buridan's Logic: The Treatise on Supposition. The Treatise on Consequences*. Synthese Historical Library n. 27. Dordrecht: D. Reidel Publishing Co. Translated, with a philosophical introduction, by Peter King.
- . 2001a. *Summulae De Demonstrationibus*. Aristarium Supplementa n. 10-8. Haren: Ingenium Publishers. Introduction, critical edition and indexes by L.M. De Rijk.
- . 2001b. *Summulae de Dialectica*. New Haven, Connecticut: Yale University Press. An annotated translation, with a philosophical introduction by Gyula Klima.

- 3722 —. 2005. *Summulae de Propositionibus*. Aristarium Supplementa n. 10-1. Nijmegen:
3723 Ingenium Publishers. Introduction, critical edition and indexes by Ria van der
3724 Lecq.
- 3725 —. 2015. *Treatise on Consequences*. New York: Fordham University Press. Translated
3726 and with an introduction by Stephen Read; editorial introduction by Hubert Hu-
3727 bien, doi:10.5422/fordham/9780823257188.001.0001.
- 3728 JOHNSTON, Spencer. 2015. “A Formal Reconstruction of Buridan’s Modal Syllogism.”
3729 *History and Philosophy of Logic* 36(1): 2–17, doi:10.1080/01445340.2014.934090.
- 3730 —. 2017. “The Modal Octagon and John Buridan’s Modal Ontology.” in *The Square of*
3731 *Opposition: A Cornerstone of Thought*, edited by Jean-Yves BÉZIAU and Gianfranco
3732 BASTI, pp. 35–52. Studies in Universal Logic. Basel: Birkhäuser, doi:10.1007/978-3-
3733 319-45062-9_4.
- 3734 KING, Peter O. 2021. “Causal Powers in the Latin Christian West.” in *Powers: A His-*
3735 *tory*, edited by Julia JORATI, pp. 112–142. Oxford Philosophical Concepts. Oxford:
3736 Oxford University Press, doi:10.1093/oso/9780190925512.003.0008.
- 3737 KUKKONEN, Taneli. 2000. “Possible Worlds in the *Tahāfut al-Falāsifa*: Al-Ghazālī on
3738 Creation and Contingency.” *Journal of the History of Philosophy* 38(4): 479–502,
3739 doi:10.1353/hph.2005.0033.
- 3740 LEFTOW, Brian. 2012. *God and Necessity*. Oxford: Oxford University Press, doi:10.109
3741 3/acprof:oso/9780199263356.001.0001.
- 3742 LEWIS, David. 1983. *Philosophical Papers, Volume 1*. Oxford: Oxford University Press,
3743 doi:10.1093/0195032047.001.0001.
- 3744 —. 1986. *On the Plurality of Worlds*. Oxford: Blackwell Publishers.
- 3745 LOUX, Michael J., ed. 1979. *The Possible and the Actual: Readings in the Metaphysics of*
3746 *Modality*. Ithaca, New York: Cornell University Press.
- 3747 NORMORE, Calvin G. 2013. “Buridanian Possibilities.” in *Logic and Language in the*
3748 *Middle Ages: A Volume in Honour of Sten Ebbesen*, edited by Jakob Leth FINK,
3749 Heine HANSEN, and Ana María MORA-MÁRQUEZ, pp. 389–402. Investigating
3750 Medieval Philosophy n. 4. Leiden: E.J. Brill, doi:10.1163/9789004242135_024.
- 3751 PAUL OF VENICE. 1990. *Logica Magna. Secunda Pars, Fasc. 4: Capitula de Condition-*
3752 *ali et de Rationali*. Toronto: Oxford University Press. Edited and translated by
3753 G.E. Hughes.
- 3754 PLANTINGA, Alvin. 1974. *The Nature of Necessity*. Oxford: Oxford University Press,
3755 doi:10.1093/0198244142.001.0001.
- 3756 PRIEST, Graham. 2016. “Thinking the Impossible.” *Philosophical Studies* 173(10):
3757 2649–2662, doi:10.1007/s11098-016-0668-5.
- 3758 SCHAFFER, Jonathan. 2010. “Monism: The Priority of the Whole.” *The Philosophical*
3759 *Review* 119(1): 31–76, doi:10.1215/00318108-2009-025.
- 3760 SCHOPENHAUER, Arthur. 1819. *Die Welt als Wille und Vorstellung*. Leipzig: Bibli-
3761 ographisches Institut, F.A. Brockhaus.

- 3762 —. 2020. *Die Welt als Wille und Vorstellung*. Hamburg: Felix Meiner Verlag. Kritische
3763 Jubiläumsausgabe der ersten Auflage von 1819 (Schopenhauer 1819) mit den
3764 Zusätzen von Arthur Schopenhauer aus seinem Handexemplar; herausgegeben
3765 von Matthias Kossler und William Massei Junior.
- 3766 SCHUMAN, Boaz Faraday. 2022. "Multiple Generality in Scholastic Logic." in *Oxford*
3767 *Studies in Medieval Philosophy*, volume X, edited by Robert PASNAU, pp. 215–282.
3768 Oxford: Oxford University Press, doi:10.1093/oso/9780192871244.003.0006.
- 3769 SHEEHY, Paul. 2006. "Theism and Modal Realism." *Religious Studies* 42(3): 315–328,
3770 doi:10.1017/s0034412506008419.
- 3771 —. 2009. "Reply to Cameron, R.P. (2009)." *Religious Studies* 45(1): 101–104, doi:10.101
3772 7/S0034412508009839.
- 3773 SOAMES, Scott. 2010. *Philosophy of Language*. Princeton, New Jersey: Princeton Uni-
3774 versity Press, doi:10.23943/princeton/9780691138664.001.0001.
- 3775 STALNAKER, Robert C. 1976. "Possible Worlds." *Noûs* 10(1): 65–75. Reprinted in Loux
3776 (1979, 225–234), in a revised form as chapter 3 of Stalnaker (1984) and in Stalnaker
3777 (2003, 26–39), doi:10.2307/2214477.
- 3778 —. 1984. *Inquiry*. Cambridge, Massachusetts: The MIT Press.
- 3779 —. 2003. *Ways a World Might Be: Metaphysical and Anti-Metaphysical Essays*. Oxford:
3780 Oxford University Press, doi:10.1093/0199251487.001.0001.
- 3781 THIJSSSEN, Hans. 2018. "Condemnation of 1277." in *The Stanford Encyclopedia of*
3782 *Philosophy*. Stanford, California: The Metaphysics Research Lab, Center for the
3783 Study of Language and Information, <https://plato.stanford.edu/archives/win2018/entries/condemnation/>.
- 3784
- 3785 THOM, Paul. 2003. *Medieval Modal Systems: Problems and Concepts*. Farnham, Surrey:
3786 Ashgate, doi:10.4324/9781315249261.
- 3787 VANCE, Chad. 2016. "Classical Theism and Modal Realism are Incompatible." *Religious*
3788 *Studies* 52(4): 561–572, doi:10.1017/S003441251600010X.
- 3789 WILLIAMSON, Timothy. 2013. *Modal Logic as Metaphysics*. Oxford: Oxford University
3790 Press, doi:10.1093/acprof:oso/9780199552078.001.0001.
- 3791 WYATT, Nicole. 2000. "Did Duns Scotus Invent Possible Worlds Semantics?" *Aus-*
3792 *tralasian Journal of Philosophy* 78(2): 196–212, doi:10.1080/00048400012349481.
- 3793 ZUPKO, John Alexander [Jack]. 2003. *John Buridan: Portrait of a Fourteenth-Century*
3794 *Arts Master*. Notre Dame, Indiana: University of Notre Dame Press.
- 3795 —. 2018. "John Buridan." in *The Stanford Encyclopedia of Philosophy*. Stanford, Cal-
3796 ifornia: The Metaphysics Research Lab, Center for the Study of Language and
3797 Information, <https://plato.stanford.edu/archives/fall2018/entries/buridan/>.

PROOF

Weakly Discerning Vertices in a Plenitude of Graphs

E. E. SHENG

De Clercq (2012) proposes a strategy for denying purported graph-theoretic counterexamples to the Principle of the Identity of Indiscernibles (PII) by assuming that any vertex is contained by multiple graphs. Duguid (2016) objects that De Clercq fails to show that the relevant vertices are discernible. Duguid is right, but De Clercq's strategy can be rescued. This note clarifies what assumptions about graph ontology are needed by De Clercq and shows that, given those assumptions, any two vertices are weakly discernible and, so, are not counterexamples to the version of PII that requires only weak discernibility.

The Principle of the Identity of Indiscernibles (hereafter PII) states that there are no *solo numero* differences. In other words, between any two things that differ numerically (i.e., differ in identity), there is a non-numerical difference (a difference that is not merely a difference in identity). Various purported counterexamples to PII have been proposed, among them Black's (1952) two intrinsically identical spheres located two miles apart in empty space. Saunders (2003) and Ladyman (2005) point out that Black's spheres and similar examples do not violate the version of PII whereby only *weak discernibility* is necessary for non-identity. A relation R weakly discerns objects a and b if and only if $Rab \ \& \ Rba \ \& \ \neg Raa \ \& \ \neg Rbb$ (Caulton and Butterfield 2012, 50). Black's spheres are weakly discerned by the relation *being two miles from*. Leitgeb and Ladyman (2008) propose cases drawn from graph theory in which, they claim, two distinct objects are not even weakly discernible. Leitgeb and Ladyman claim that—whereas the two vertices in the graph consisting of two vertices and an edge connecting them are weakly discernible—the two vertices in the graph consisting of two vertices and no edges are not in any way discernible. De Clercq (2012) argues that Leitgeb and Ladyman's counterexample rests on a controversial view about the ontology of graphs, namely, one that rejects

3828 assumption (i) below, and that on another plausible view about the ontology of
 3829 graphs, which De Clercq favours, the case that Leitgeb and Ladyman propose
 3830 is not a counterexample to PII, because the two vertices are discernible in
 3831 virtue of the relations in which they stand in other graphs that contain them.
 3832 Duguid (2016) objects that the two vertices are not discernible in virtue of
 3833 such relations, so that, even granting De Clercq's favoured view about the
 3834 ontology of graphs, Leitgeb and Ladyman's case is a counterexample to PII,
 3835 even the version of PII that requires only weak discernibility.

3836 In this note, I clarify what assumptions about the ontology of graphs are
 3837 needed by De Clercq and show that De Clercq's strategy can be rescued from
 3838 Duguid's rejoinder insofar as it can be shown that, granted De Clercq's assump-
 3839 tions about the ontology of graphs, any two vertices are weakly discernible. I
 3840 give an example of a relation that weakly discerns vertices: *x has greater degree*
 3841 *in some graph than y*. If De Clercq is correct about the ontology of graphs,
 3842 therefore, the purported graph-theoretic counterexamples that have been pro-
 3843 posed do not falsify the version of PII that requires only weak discernibility
 3844 and thus do not, in this respect, improve on Black's spheres.

3845 1 De Clercq's Strategy

3846 Graphs are arrangements of vertices and edges connecting vertices such that
 3847 edges do not have a direction and any two vertices in a graph are either
 3848 connected by one edge or not connected by any edge.¹ The *degree* of a vertex
 3849 in a graph is the number of edges that connect it with other vertices in the
 3850 graph. A vertex is *isolated* in a graph if and only if it has degree 0 in the graph.
 3851 Two graphs are *isomorphic* if and only if (regardless of the identities of their
 3852 vertices) they have the same structure of vertices and edges; that is, two graphs
 3853 are isomorphic if and only if there is a bijection from the set of the vertices of
 3854 the first graph to the set of the vertices of the second graph such that any two
 3855 vertices are connected by an edge in the first graph if and only if their images
 3856 under the bijection are connected by an edge in the second graph.

3857 More formally, graphs are commonly defined set-theoretically so that a
 3858 *graph* G is an ordered pair (V, E) where V is a set of *vertices* and E is a (possibly
 3859 empty) set of subsets of V that have two members, and any distinct vertices v
 3860 and w in V are said to be connected in (V, E) by an *edge* if and only if $\{v, w\}$

1 In *directed graphs*, edges have a direction. In *multigraphs*, vertices may be connected by more than one edge. Directed graphs and multigraphs are not considered in this note.

3861 is a member of E .² Let us call the identification of graphs with ordered pairs
 3862 of vertex and edge sets *Identity*. Leitgeb and Ladyman do not accept *Identity*,
 3863 while De Clercq does. Other graph-theoretic terms can also be defined set-
 3864 theoretically.

3865 De Clercq defends PII against Leitgeb and Ladyman's purported counterex-
 3866 ample by arguing that, in a graph G_0 that consists of two vertices, a and b ,
 3867 and no edges, a and b are discernible in virtue of the relations in which they
 3868 stand in other graphs: "vertices in labeled graphs are always distinguishable,
 3869 not just because they bear different labels, but also because they feature in
 3870 (structurally!) different ways in different graphs" (2012, 670). The distinct-
 3871 ness of a and b , for example, is, according to De Clercq, not a *solo numero*
 3872 difference, because there is a graph G_2 , consisting of three vertices a , b , and c
 3873 (where a and b are respectively identical to the vertices a and b in G_0) and an
 3874 edge connecting a and c , in which a and b stand in different relations.

3875 Two assumptions are necessary and sufficient for De Clercq's strategy: (i)
 3876 that there are no *unlabelled graphs* such that their vertices are objects, and
 3877 (ii) that if G_0 exists, then G_2 exists. Regarding (i): De Clercq and Leitgeb and
 3878 Ladyman disagree about what unlabelled graphs are. According to Leitgeb and
 3879 Ladyman, unlabelled graphs are graphs such that the vertices of an unlabelled
 3880 graph are distinguished only by their relations within that unlabelled graph,
 3881 and any isomorphic unlabelled graphs are identical. On this view, there are
 3882 objects that are the vertices of unlabelled graphs, and vertices of distinct
 3883 unlabelled graphs do not stand in relations of identity. De Clercq (2012, 666),
 3884 in contrast, claims that "unlabeled graphs are not graphs but isomorphism
 3885 classes of graphs" (that is, the equivalence classes into which the set of all
 3886 graphs is partitioned by the isomorphism relation).³ On this view, talk of
 3887 the vertices of unlabelled graphs is not ontologically committing, and there
 3888 are no unlabelled graphs such that their vertices are objects. Regarding (ii):
 3889 Since, as specified above, G_2 is a graph some of whose vertices are respectively
 3890 identical to some vertices of G_0 , (ii) presupposes (iii) that some vertices in

2 Note: This definition is not committed to identifying edges with sets of two vertices.

3 Note: De Clercq's identification of unlabelled graphs with isomorphism classes is not necessary for his argument. One could instead claim, for example, that unlabelled graphs are mereological atoms that correspond one-to-one with isomorphism classes. But perhaps it is the best motivated of claims that imply (i).

3891 distinct graphs are identical.⁴ De Clercq (2012, 665–669) defends (i) and (iii)
 3892 by appealing to the practice of graph theorists.⁵

3893 Assuming uncontroversially that there exist graphs of three or more vertices
 3894 and thus that there exists at least one vertex other than a and b , (ii) follows
 3895 from the following claim:

3896 **PLENITUDE.** For every subset V of the set W of all vertices, for every
 3897 (possibly empty) set E of sets consisting of two members of V , there
 3898 exists a graph consisting of the vertices in V and edges connecting
 3899 every pair of members u, v of V such that $\{u, v\}$ is in E .⁶

3900 In turn, **PLENITUDE** follows from *Identity*, since any set exists if its members
 3901 exist, and any ordered pair exists if its members exist. So, De Clercq can accept
 3902 *Identity* and infer (ii) from *Identity*, but only if he also accepts **PLENITUDE**.
 3903 Might one accept (ii) without accepting **PLENITUDE**? As noted above, (ii)
 3904 implies (iii) that some vertices in distinct graphs are identical, or, in other
 3905 words, that some vertices are contained by multiple graphs. De Clercq (2012,
 3906 666–667) argues that, while (iii) follows from *Identity*, it is also plausible in
 3907 light of mathematical practice, independently of the truth of *Identity*. None-
 3908 theless, as long as some vertices are contained by multiple graphs, it would
 3909 be arbitrary to suppose that some finite graphs that can be formed out of ver-
 3910 tices from W and edges connecting them exist but others do not. De Clercq’s
 3911 assumption of (ii), therefore, commits him to **PLENITUDE**.

3912 It is **PLENITUDE** that leaves De Clercq’s defence of PII vulnerable, even
 3913 granting (i), to Duguid’s reply: for any graph where a and b bear different
 3914 relations, another graph exists in which a and b are permuted so that (for
 3915 instance) corresponding to G_2 there exists an isomorphic graph G_1 that consists
 3916 of three vertices a, b , and c , and an edge connecting b and c . Now, b has the
 3917 property *being isolated in a graph consisting of three vertices and an edge*
 3918 *connecting two of them* in virtue of G_2 , but a has the same property in virtue

4 Note: *Identity* is not necessary for De Clercq’s strategy because (i) and (ii) are sufficient for it and do not imply *Identity*. *Identity* is also not sufficient for De Clercq’s strategy because, although *Identity* implies that there are no graphs that are Leitgeb and Ladyman’s unlabelled graphs, *Identity* does not imply that Leitgeb and Ladyman’s unlabelled graphs do not exist, and as long as such entities exist, there are counterexamples to PII.

5 In defending the rejection of (i), Leitgeb and Ladyman (2008, 390) also appeal to the practice of graph theorists.

6 Note: **PLENITUDE**, thus formulated, does not presuppose *Identity*, as it would if “a graph consisting [...] is in E ” were replaced with “a graph (V, E) .”

3919 of G_1 . To discern a and b , De Clercq would have to appeal to properties
 3920 that distinguish between isomorphic graphs (for example, the property *being*
 3921 *isolated in G_2*).⁷ But, since distinct isomorphic graphs differ only in the identity
 3922 of their vertices, in order to distinguish between isomorphic graphs, a property
 3923 “must utilize object names” (Duguid 2016, 472). Let us say that a property or
 3924 relation is *forbidden* for PII if and only if allowing being discerned by it to count
 3925 as discernibility would make PII metaphysically uninteresting. For example, a
 3926 version of PII that allows that a and b are discernible on the ground that they
 3927 are discerned by the property *is identical to a* is metaphysically uninteresting,
 3928 as is a version of PII that allows that a and b are discernible on the ground that
 3929 they are weakly discerned by the relation *is distinct from*. (Rodriguez-Pereyra
 3930 2006; and Muller 2015 discuss what would make PII trivial and, as such,
 3931 metaphysically uninteresting.) Following Muller (2015), Duguid considers
 3932 properties in which object names occur to be forbidden for PII. De Clercq,
 3933 Duguid concludes, fails to save PII.

3932 Weakly Discerning Vertices

3935 Ladyman, Linnebo and Pettigrew (2012) show in their theorem 6.4 that two
 3936 objects are weakly discernible in a language L if and only if they are in any
 3937 way discernible in the language that includes a constant for every element of
 3938 the domain of L (i.e., the language that includes names for all of its objects). It
 3939 follows, as Duguid accepts, that if there are object-name-containing properties
 3940 that discern two vertices, a and b , there is a non-object-name-containing
 3941 relation that weakly discerns a and b . Nonetheless, Duguid (2016, 473) writes,
 3942 “Such a relation has not yet been provided. And neither can I see what it might
 3943 be.”

3944 Here is one:

$$\Phi(x, y) := \exists g((g \text{ is a graph}) \ \& \ (g \text{ contains } x \text{ and } y) \ \& \\ (x \text{ has greater degree in } g \text{ than } y))$$

3945 Given De Clercq’s assumptions, this relation, *x has greater degree in some graph*
 3946 *than y*, holds between any two vertices, a and b , in both directions, but not

7 Duguid (2016, 472) says that De Clercq must appeal to a property that is “specific enough to single out a single graph.” This is not correct, since the property *is isolated in a graph consisting of a , b , and some third vertex and an edge connecting a and the third vertex*, which does not single out a single graph, would also do.

between either vertex and itself. Hence, *contra* Duguid, any two vertices are weakly discernible, and Leitgeb and Ladyman's case is not a counterexample to the version of PII that requires only weak discernibility.

Whether De Clercq's strategy for saving PII from purported graph-theoretic counterexamples is ultimately successful depends on the plausibility of its assumptions about graph ontology: **PLENITUDE** and that there are no unlabelled graphs such that their vertices are objects. Granted these assumptions, however, any two vertices are indeed weakly discernible.*

E. E. Sheng

 0000-0001-9324-7704

University of Oxford

eric.sheng@philosophy.ox.ac.uk

References

- BLACK, Max. 1952. "The Identity of Indiscernibles." *Mind* 61(242): 153–164. Reprinted in Black (1954, 80–92), doi:10.1093/mind/LXI.242.153.
- . 1954. *Problems of Analysis: Philosophical Essays*. Ithaca, New York: Cornell University Press.
- CAULTON, Adam and BUTTERFIELD, Jeremy. 2012. "On Kinds of Indiscernibility in Logic and Metaphysics." *The British Journal for the Philosophy of Science* 63(1): 27–84, doi:10.1093/bjps/axr007.
- DE CLERCQ, Rafael. 2012. "On Some Putative Graph-Theoretic Counterexamples to the Principle of the Identity of Indiscernibles." *Synthese* 187(2): 661–672, doi:10.1007/s11229-010-9867-3.
- DUGUID, Callum. 2016. "Graph Theory and the Identity of Indiscernibles." *Dialectica* 70(3): 463–474, doi:10.1111/1746-8361.12151.
- LADYMAN, James. 2005. "Mathematical Structuralism and the Identity of Indiscernibles." *Analysis* 65(3): 218–221, doi:10.1111/j.1467-8284.2005.00552.x.
- LADYMAN, James, LINNEBO, Øystein and PETTIGREW, Richard. 2012. "Identity and Discernibility in Philosophy and Logic." *The Review of Symbolic Logic* 5(1): 162–186, doi:10.1017/S1755020311000281.
- LEITGEB, Hannes and LADYMAN, James. 2008. "Criteria of Identity and Structuralist Ontology." *Philosophia Mathematica* 16(3): 388–396, doi:10.1093/phimat/nkm039.
- MULLER, F. A. 2015. "The Rise of Relationals." *Mind* 124(493): 201–237, doi:10.1093/mind/fzu175.
- RODRIGUEZ-PEREYRA, Gonzalo. 2006. "How Not to Trivialize the Identity of Indiscernibles." in *Universals, Concepts and Qualities: New Essays on the Meaning of*

* THANKS

- 3983 *Predicates*, edited by Peter Frederick STRAWSON and Arindam CHAKRABARTI, pp.
3984 205–223. Farnham, Surrey: Ashgate, doi:[10.4324/9781315235325-17](https://doi.org/10.4324/9781315235325-17).
3985 SAUNDERS, Simon W. 2003. “Physics and Leibniz’s Principles.” in *Symmetries in Physics:*
3986 *Philosophical Reflections*, edited by Katherine BRADING and Elena CASTELLANI,
3987 pp. 289–307. Cambridge: Cambridge University Press, doi:[10.1017/cbo978051153](https://doi.org/10.1017/cbo9780511535369.017)
3988 [5369.017](https://doi.org/10.1017/cbo9780511535369.017).

PROOF

PROOF

3989

Published by *Philosophie.ch*

3990

Verein philosophie.ch

3991

Fabrikgässli 1

3992

2502 Biel/Bienne

3993

Switzerland

3994

dialectica@philosophie.ch

3995

<https://dialectica.philosophie.ch/>

3996

ISSN 0012-2017

3997

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

3998

3999

Dialectica is supported by the [Swiss Academy of Humanities and Social Sciences](https://www.snf.ch/).

4000

4001

Abstracting and Indexing Services

4002

The journal is indexed by the Arts and Humanities Citation Index, Current

4003

Contents, Current Mathematical Publications, Dietrich's Index

4004

Philosophicus, IBZ — Internationale Bibliographie der Geistes- und

4005

Sozialwissenschaftlichen Zeitschriftenliteratur, Internationale Bibliographie

4006

der Rezensionen Geistes- und Sozialwissenschaftlicher Literatur, Linguistics

4007

and Language Behavior Abstracts, Mathematical Reviews, MathSciNet,

4008

Periodicals Contents Index, Philosopher's Index, Repertoire Bibliographique

4009

de la Philosophie, Russian Academy of Sciences Bibliographies.

4010 **Contents**

4011 LÁSZLÓ BERNÁTH & TAMÁS PAÁR, *Responsibility First: How to Resist Ag-*
4012 *nosticism About Moral Responsibility* 525

4013 BENJAMIN WINOKUR, *Self-Knowledge and Interpersonal Reasoning* 547

4014 PAUL O. IRIKEFE, *Avner Baz’s Ordinary Language Challenge to the Philo-*
4015 *sophical Method of Cases* 571

4016 DAN MARSHALL, *The Dual-Detector Argument Against the Modal Theory* . . 597

4017 BOAZ FARADAY SCHUMAN, *Lewisian Worlds and Buridanian Possibilia* 623

4018 E. E. SHENG, *Weakly Discerning Vertices in a Plenitude of Graphs* 651

PROOF