

# dialectica

International Journal of Philosophy

## Contents

GRAHAM STEVENS, <i>Metalinguistic Monstrosity and Displaced Communications</i> .....	615
RYAN COX, <i>Constitutivism about Instrumental Desire and Introspective Belief</i>	653
ALIDA LIBERMAN, <i>The Mental States First Theory of Promising</i> .....	683
PETER FINOCCHIARO, <i>A Puzzle About Parsimony</i> .....	711
TOM DONALDSON, <i>David Armstrong on the Metaphysics of Mathematics</i> ....	729

PROOF

# 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 (also managing editor)

## Review Editors

Stephan Leuenberger and Philipp Blum

## Editorial Committee

Philipp Blum (né Keller), Claudio Calosi, Zoé Christoff, Fabrice Correia, Catharine Diehl, Matthias Egg, Patrik Engisch, Andrea Giananti, Jörg Löschke, Arturs Logins, Giovanni Merlo, Robert Michels, Ryan Miller, Paolo Natali, Donnchadh O'Conaill, François Pellet, Edgar Phillips, Stephanie Rennick, Maria Scarpati, Mike Stuart, Fabrice Teroni, Daniel Vanello, Lisa Vogt.

## 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), Thomas Strahm (Bern), Christine Tappolet (Montréal), Neil Tennant (Ohio State), Mark Textor (King's College London), Achille Varzi (Columbia University), Alberto Voltolini (Torino), Timothy Williamson (Oxford).

December 2020

## Contents

GRAHAM STEVENS, <i>Metalinguistic Monstrosity and Displaced Communications</i> .....	615
RYAN COX, <i>Constitutivism about Instrumental Desire and Introspective Belief</i>	653
ALIDA LIBERMAN, <i>The Mental States First Theory of Promising</i> .....	683
PETER FINOCCHIARO, <i>A Puzzle About Parsimony</i> .....	711
TOM DONALDSON, <i>David Armstrong on the Metaphysics of Mathematics</i> ....	729

PROOF

# Metalinguistic Monstrosity and Displaced Communications

GRAHAM STEVENS

David Kaplan's semantic theory for indexicals yields a distinct logic for indexical languages that generates contingent a priori truths. These special truths of the logic of indexicals include examples like "I am here now," an utterance of which expresses a contingent state of affairs and yet which, according to Kaplan, cannot fail to be true when it is uttered. This claim is threatened by the problem of *displaced communications*: answerphone messages, for example, seem to facilitate true instances of the negation of this supposed logical truth as they allow the agent of the message to no longer be at the location of the message when it is encountered by an audience. Many such displaced communications can be identified in everyday natural language uses of indexicals. Recent discussion has suggested that Kaplan's error is to be overly restrictive in the possible contexts of utterance his semantic theory recognizes, as he fails to acknowledge the possibility of utterances that occur at a context distinct from that in which they are constructed. I reject this diagnosis and defend Kaplan's semantic theory. Displaced communications, I argue, are best understood as resulting from a pragmatically introduced metalinguistic context-shifting operation and hence do not demand revision of Kaplan's semantic theory. I provide an analysis of the pragmatic process underlying this operation and make the case for its merits over those of rival accounts of displaced communications.

David Kaplan's (1989a) semantic theory for indexicals yields a distinctive logic for indexical languages, generating a set of logical truths that are entirely absent from non-indexical languages. These logical truths are notable in that they invalidate the rule of *necessitation* ( $\vDash \phi \rightarrow \vDash \Box\phi$ ), because there are sentences that Kaplan thinks cannot be uttered without being true, despite the fact that they express non-necessary states of affairs. Consequently, they are often cited as examples of contingent a priori truths. Recent discussion of indexicality in the philosophical literature has challenged Kaplan's

32 proposal to grant this privileged status to certain indexical constructions,  
33 however, by drawing attention to numerous apparent counter-examples in  
34 natural language. These challenges almost unanimously agree that Kaplan is  
35 too restrictive in his analysis of the sorts of contexts in which an indexical  
36 sentence can be employed.<sup>1</sup> All of the proposed counter-examples appear to  
37 show that under certain conditions uses of indexical sentences can align an  
38 indexical sentence with a context that is not recognized by Kaplan's theory  
39 and, therefore, that Kaplan's apparent cases of contingent a priori truths do  
40 not reflect genuine semantic features of English, but only reflect Kaplan's  
41 mistaken intuitions about the admissible range of contexts in which indexical  
42 sentences may be uttered.

43 In this paper I will defend Kaplan's semantic theory against this challenge.  
44 I will proceed by first arguing that the proposed counter-examples in ques-  
45 tion are not just the result of aligning an indexical sentence with an unusual  
46 context, they are the result of applying a context-shifting operator on the  
47 character of an indexical sentence. Kaplan calls an operator on character a  
48 "monster" and argues that monsters are entirely absent from the semantics of  
49 English. There are, however, metalinguistic devices such as quotation that do  
50 behave like monsters, as acknowledged by Kaplan. I will argue that the pro-  
51 posed counter-examples to Kaplan's account all share important similarities  
52 with these metalinguistic operators and are thus best understood as resulting  
53 from pragmatically introduced metalinguistic operators on constructions that  
54 are perfectly acceptable on Kaplan's analysis. I thus conclude that Kaplan's  
55 semantic theory does not stand in need of revision to accommodate these  
56 examples, and we have every reason to retain the view that indexicals can gen-  
57 erate the sorts of contingent a priori truths predicted by Kaplan. I will begin  
58 by elucidating Kaplan's theory, then discussing the challenges to that theory. I  
59 will next introduce the notion of a monster and present the argument that all  
60 of the proposed challenges depend on what I will call *monstrous operations*. I  
61 will then defend the view that these particular monstrous operations are best  
62 explained pragmatically as resulting from metalinguistic operations, rather  
63 than semantic ones. Finally, I will consider some objections and replies.

---

1 Some notable exceptions are discussed below.

## 6.1 The Logic of Indexicals

65 Pure indexicals are expressions whose literal meaning both requires a context  
66 for saturation and specifies precisely what role context must play in the satu-  
67 ration (in English, examples include “I,” “now,” “today,” some uses of “here,”  
68 etc.). Demonstratives differ from pure indexicals insofar as they require an  
69 additional demonstration such as a gesture or other such directing intention  
70 (in English, examples include “that,” “this,” “she,” some other uses of “here,”  
71 etc.). To account for this distinctive class of meanings, Kaplan proposes a two-  
72 level semantic theory, coupled with a double-indexing of indexical sentences  
73 with formal representations of contextual situations. Firstly, an indexical ex-  
74 pression is associated with both a *character* and a *content*. The character of  
75 the expression is a function from *context* to content. Less formally expressed,  
76 this means that the character can be thought of as a rule governing the con-  
77 tribution required by context in order to fix the semantic value or reference of  
78 the term. For example, the first-person English pronoun “I” has a character  
79 of the form “the agent of this utterance.” This description specifies a function  
80 that will yield a different value depending on who is speaking. Sentences  
81 have the same two-level semantic profile and each level is compositionally  
82 derived from the individual expressions contained in the sentence. For exam-  
83 ple, the sentence “I am walking” has a character that maps a context  $c$   
84 on to the proposition that  $a_c$  is walking, where  $a_c$  is the agent of the context  
85  $c$ . Different contexts will be mapped to different contents by this function.  
86 If character is a function from contexts to contents, then individual agents,  
87 objects, times, places and so on will be the contents of individual indexical  
88 and demonstrative expressions when they are used in context. The content  
89 of an indexical is thus its reference.<sup>2</sup> The content of a sentence is a propo-  
90 sition. It follows that indexical sentences cannot contribute truth-evaluable  
91 propositions to a semantic theory without the assistance of context. Hence,  
92 within the semantic theory, it is ordered pairs of sentences and contexts that

---

2 An insistence that *all* indexicals *must* be directly referential is misplaced if we take the definition of an indexical simply to be that its character is a non-constant function from contexts to contents. There seems no principled reason to exclude quantifiers, predicates, or unarticulated constituents such as those which are regularly posited to fix the comparison class for gradable adjectives, from this list. King’s (2001) quantificational analysis of complex demonstratives, for example, treats such expressions as indexical quantifiers. Similarly Kaplan’s formal language *LD* contains an indexical temporal operator in place of a referring expression as its correlate of “now.” Nonetheless it seems obvious that most paradigm cases of indexicals and demonstratives (“I,” “here,” “this,” etc.) are directly referential.

are the bearers of truth-values. Contexts themselves are precisely specified as sequences of parameters corresponding to the indexical elements in the sentence. So, for example, a sentence containing “I,” “here,” and “now” will demand a context with parameters for an agent, a spatial location, and a temporal location. We also add a world parameter to every context. Thus the context in this case will be of the form  $c = \langle a_c, l_c, t_c, w_c \rangle$ . The sentence-context pair  $\langle s, c \rangle$  models the proposition expressed by the sentence  $s$  in the context  $c$ . Utterances of sentences are thus indexed to contexts. Similarly, evaluation of propositions requires what Kaplan calls “circumstances of evaluation,” which are (at a minimum) pairs of world and times. As we will now see, this double-indexing reveals that some sentences will be true with respect to any context they are paired with, despite expressing propositions which are not true at every circumstance of evaluation (thus, are not necessary).

With the above system outlined, we can make sense of Kaplan’s claim to have discovered examples of the contingent a priori in English. Consider the following four sentences:

- (a) I am me.
- (b) This is of the same chemical kind as that.
- (c) I am here now.
- (d) I am not alive.

Any utterance of (a) will be true because any context  $c$  will be such that  $a_c = a_c$ . Thus we know a priori that any utterance of (a) will be true. And, as this truth holds in all possible worlds, (a) will be true at every circumstance. Thus (a) is a necessary truth. An utterance of (b) in a context  $c_1$  in which the demonstratum of “this” is a sample of liquid water, and the demonstratum of “that” is a sample of water-ice, will be true. But an utterance of the same sentence with different demonstrata could well be false. Hence it is certainly not true a priori. However, if it is uttered in  $c_1$ , then the truth it expresses will hold at every circumstance of evaluation. Thus it is an example of a necessary a posteriori truth. Kaplan also holds that any utterance of (c) must be true, on the grounds that there is no context in which an agent can fail to be at the location of that context at the time of that context when uttering something in that context. All the same, their being located at whatever part of space-time they are at when they make that utterance is obviously just a contingent fact—they could have been elsewhere. Hence the proposition expressed does not hold true at every circumstance. Accordingly (c) is assumed by Kaplan to



129 be a contingent a priori truth. Likewise for (d), no-one can utter this truly (as  
130 the saying goes, dead men tell no tales), but there is nothing necessary about  
131 one's being alive. Thus an utterance of (d) is known to be false on a priori  
132 grounds, but is not *necessarily* false.

133 What makes Kaplan's defence of the contingent a priori so compelling is that  
134 Kaplan's proposed cases require no investment in any kind of metaphysical  
135 speculation. They are just immediate consequences of the correct semantic  
136 analysis of indexicals. Or so it seemed. Many have reacted to Kaplan's logic  
137 of indexicals, however, by pointing out that Kaplan's analysis does not seem  
138 to be proceeding on purely semantic grounds but in fact makes significant  
139 assumptions about the conditions under which utterances and other forms  
140 of communication can be made that are empirically questionable. Consider  
141 again the example (c) above. This apparent logical truth is routinely negated  
142 as an answerphone recording: "I am not here now. Please leave a message  
143 after the tone...". Similarly, one may record a message to be replayed at the  
144 reading of one's will that contains (d): "If you are hearing this recording, then  
145 I am not alive. I have left you this message to communicate my wishes to you  
146 after my death...". In such circumstances it appears that truths, not logical  
147 falsehoods as Kaplan's analysis seems to predict, are being communicated.

148 These counter-examples to Kaplan's analysis highlight, and challenge, an  
149 assumption in Kaplan's theory about the interplay between sentence and  
150 context. Kaplan's assumption is that contexts of utterance, inscription, or  
151 other sorts of linguistic performance, always conform to a minimum structural  
152 norm such that agents of utterances are always located in the time and place  
153 of that utterance. Indeed Kaplan explicitly acknowledges this, arguing that  
154 we must restrict contexts of utterance to these "proper" contexts: "[I]mproper  
155 indices are like impossible worlds; no such contexts could exist and thus there  
156 is no interest in evaluating the extensions of expressions with respect to them"  
157 (1989a, 509). Prima facie, these counter-examples seem to show that Kaplan  
158 is wrong. Furthermore, the counter-examples are ubiquitous. Here is another,  
159 from Predelli (2005, 43). Jones writes the following note just before leaving  
160 his house at 8am, which he then leaves at home for his wife, who is not due  
161 to return until 5pm:

162 I am not at home now. If you hurry, you'll catch the evening flight  
163 to Los Cabos. Meet me in six hours at the Hotel Cabo Real.

164 In this example, the note is obviously not intended to be, nor will it be, inter-  
 165 preted as indexed to the time at which it was inscribed but, rather, is intended  
 166 to be indexed to the time at which it will be read. Another example, offered  
 167 by Corazza, Fish and Gorvett (2002), challenges Kaplan's assumption that the  
 168 agent, utterer, and referent, of "I" in a given context must always coincide.  
 169 They invite us to consider the case of an academic who leaves a note on his  
 170 door saying "I am not here today" to signal his absence when it is read. This,  
 171 already, is an example equivalent to Predelli's above, but they continue the  
 172 story by imagining that the academic returns to work several days later and  
 173 then removes the note and reattaches it to a colleague's door to signal *their*  
 174 absence. It now seems that the referent of "I" (along with other indexicals)  
 175 has shifted while its inscriber and possibly even agent have not. What are we  
 176 to say of these challenges to Kaplan's theory? In the next section I will argue  
 177 that all of these counter-examples display "monstrous" properties.

## 172 **Monsters and Monstrous Operations**

179 Kaplan gives the name "monster" to any would-be operator on the character  
 180 of an indexical. He maintains that no such operators exist in English. Take for  
 181 example, the indexical "I." A monster operating on this expression would effect  
 182 a context shift such that the reference of the expression shifted on to someone  
 183 other than the agent of the context of utterance. But no such operation seems  
 184 possible within the confines of ordinary English. If I say "in some contexts, I  
 185 am not me," this is simply false (if interpreted literally—there may of course be  
 186 figurative uses of this sentence which are understood to communicate a truth).  
 187 I am always identical to myself. Similarly, if I embed the same indexical within  
 188 a propositional attitude verb, the attitude verb has no impact on the character  
 189 of the indexical, which immediately takes wide scope with respect to it: an  
 190 utterance of "John believes that I am happy" communicates the speaker's  
 191 report that John believes the speaker to be happy. The character of "I" picks  
 192 out the speaker, regardless of any operators contained in the utterance.

193 Kaplan does point out, however, that monstrous operations can be cre-  
 194 ated by *metalinguistic* devices. The most common is quotation. Compare the  
 195 following:

- 196 (e) John said that I am happy.
- 197 (f) John said "I am happy."

198 By naming the indexical sentence “I am happy” we can shift the character  
199 of that sentence in (f), whereas it is impossible in (e). Kaplan’s position,  
200 then, is that the only operators resembling monsters that can be applied  
201 to English expressions are metalinguistic operators. To keep this distinction  
202 between linguistic and metalinguistic operations intact in what follows, I  
203 will use the expression “monster” to denote the sort of lexicalized linguistic  
204 operator that Kaplan maintains is absent from English,<sup>3</sup> and the expression  
205 “monstrous operator” (hereafter “MO”) to denote any operation, including  
206 the metalinguistic operators such as quotation, on character. Thus, according  
207 to my usage, every monster is an MO but not every MO is a monster.

### 208 3 Displaced Communications

209 The counter-examples to Kaplan’s proposed truths of the logic of indexicals  
210 have been responded to in varying ways, but most of the responses conform  
211 to one general strategy. The counter-examples are usually understood as  
212 somehow involving a displacement from its point of origin of the information  
213 that is communicated. To put it another way, a distinction is drawn between  
214 the point at which the communication is encoded and the point at which  
215 it is decoded. The alleged flaw in Kaplan’s reasoning has thus been almost  
216 unanimously identified as the mistaken assumption that communications  
217 occur at the point when (and where) they are encoded (either recorded in the  
218 cases like answerphone messages, or inscribed in the case of written notes  
219 and messages). By contrast, these counter-examples all seem to be intuitively  
220 understood as communicating information at the time when they are decoded.  
221 Sidelle (1991, 535) describes the production of an answerphone message as  
222 a process of “arranging to make an utterance at a later time, or, if one likes,  
223 deferring an utterance.” This interpretation has gone largely unchallenged  
224 despite the differences in opinion as to the semantic or pragmatic mechanism  
225 by which this proposed procedure is thought to be realized.<sup>4</sup> In what follows

---

3 Kaplan only says that monsters are absent from English but he is often interpreted as making the wider claim that they are absent from natural languages generally, a claim challenged by Schlenker (2003), who appeals to empirical data concerning Amharic to support the view that monsters are present in some natural languages. For a detailed discussion of Kaplan on monsters, see Predelli (2014).

4 It is denied by Stevens (2009). Cohen (2013, 8, fn 8) rightly points out that a major shortcoming of Stevens (2009) is my lack of a positive proposal in place of this interpretation. In particular, no detail is given as to how a pragmatic explanation could explain apparent utterances at a distance, so as to make the idea of a deferred utterance redundant. A proposal like the one I will be offering

226 I will use the term “Displaced Communication” (hereafter “DC”) to denote  
 227 this proposed act of encoding content for later decoding.<sup>5</sup>

228 It is rarely, if indeed ever, noted that any such act of encoding content for  
 229 later communication is itself an MO in the sense that all such cases involve  
 230 the use of an indexical to communicate a content which is shifted away from  
 231 the context in which the communication is encoded. For example, the use of  
 232 “now” to refer to a time other than the time at which it is being used to encode  
 233 the communication, the use of “I” to refer to someone other than the agent  
 234 who is encoding the communication, and so on. Furthermore, it is clear that  
 235 the operation of deferring the communication must involve this monstrous  
 236 operation. If it did not, then we would not have our counter-examples to  
 237 Kaplan’s proposed analysis. DCs are not monsters in the sense of being discrete  
 238 lexically encoded operators within the language that introduce MOs, but they  
 239 are certainly monstrous in this wider sense that includes pragmatic processes  
 240 and operations (which is clearly encompassed by the definition of an MO  
 241 provided above).

242 Two objections may be raised against the above description of displaced  
 243 communications as being (or being the result of) MOs, which should be ad-  
 244 dressed before continuing. Firstly, one may point to the fact that displaced  
 245 communications appear to perform (or result from the performance of) an  
 246 operation directly on the *context* with which the communicated sentence is  
 247 paired, rather than the *character* of the sentence. Secondly, one may point  
 248 to the fact that monsters, as traditionally understood, play a recognisably se-  
 249 mantic role by operating directly on lexical items through binding operations,  
 250 imposing scope relations on them, and so forth. I will reply to each objection  
 251 in turn.

252 On Predelli’s (2005) interpretation of DCs, the correct context that a sen-  
 253 tence must be paired with for accurate semantic evaluation is determined by  
 254 the intention of the speaker. That I can intend my utterance at time  $t_1$  of the  
 255 sentence “I am not here now” to be evaluated at a later time  $t_2$  demonstrates,  
 256 on this view, Kaplan’s error in defining proper contexts too narrowly. On

---

in this paper is going to be required if we are to reject the deferred utterance analysis. Others have offered pragmatic proposals that share some features with my approach, including Connolly (2017) who also diagnoses the counterexamples as involving some form of pretence, although his analysis of how this is effected takes a very different line to mine. Åkerman (2017) also provides an alternative account of how pragmatic processes can be appealed to in our explanation of apparent cases of context-shifting.

5 I prefer the term “displaced communication” to “deferred utterance” as it is not limited to utterances, but can include any communication of information.

257 this view, the counter-examples to Kaplan's theory arise because sentences  
258 have been paired with contexts that he did not recognize. As such, it does not  
259 automatically appear that any MO has been applied. The character, Predelli  
260 (2005, 44) insists, does not change at all. Rather, careful inspection of the  
261 "preparatory operations" (2005, 58)—a phrase Predelli borrows from Quine to  
262 describe the decisions we make about how to regiment particular utterances  
263 or inscriptions to make them ready for semantic evaluation—simply reveals  
264 that sentences can in fact be paired with a wider selection of contexts than  
265 Kaplan anticipated. All of the work that separates Predelli's position from  
266 Kaplan's occurs at this pre-semantic stage: "[O]nce the appropriate clause-  
267 index pair has been identified, the indexicals proceed with their customary  
268 characters, and results of truth-value are obtained on the basis of the usual  
269 mechanisms of compositional analysis" (2005, 58). Nonetheless, further re-  
270 flection makes it clear enough that this position entails the existence of MOs.  
271 An MO is an operation which shifts the context of an indexical element or  
272 sentence by operating on its character. We can easily modify the examples of  
273 displaced communication to make this effect more explicit. For example, one  
274 can download an audio file of Arnold Schwarzenegger's famous utterance  
275 of "I'll be back" from the *Terminator* movie and set this file to play as an an-  
276 swerphone message. This performs an MO on the original utterance (made by  
277 Schwarzenegger), shifting the context to one in which the agent is no longer  
278 Arnold's cyborg character from the movie, but a real person who has set up  
279 an answerphone message to denote their temporary absence or unavailability  
280 to a caller.

281 A further worry may arise by comparing these sorts of MOs with the oper-  
282 ators that Kaplan defines as monsters. A monster is an operator on the  
283 character of an expression. One way to identify the presence of a monster  
284 would be to look for the observable effects that the monster has on the scope  
285 of the expression operated on. For example, if we try to shift the scope of  
286 the indexical "now" by using a phrase like "in some contexts," we get the  
287 following:

288 (g) In some contexts, yesterday is now.

289 The attempt to shift the context fails because "now" and "yesterday" resist  
290 embedding under the scope of the operator; i.e. they take wide scope over the  
291 operator. On the surface, this seems to be a quite different operation to any  
292 present in cases of DCs. The answerphone message "I am not here now," for

293 example, does not display any distinctive impact on the scope of the indexicals  
 294 contained in it. While this is undeniable, I don't think it counts against the  
 295 view that an **MO** is at work in cases of **DCs**.

296 For one thing, there are many examples of **DCs** which are quite naturally  
 297 understood as imposing a binding operation on an indexical element. Take for  
 298 example the logo used on merchandise by the Rock Climbing equipment man-  
 299 ufacturer DMM. They produce clothing with the following phrase emblazoned  
 300 on it:

301 (h) Climb now, work later.

302 This phrase can be naturally used in conversation in a way that pairs it with a  
 303 proper Kaplanian context:

304 *Speaker A:* I don't know whether I should go climbing now, or after  
 305 I have finished my essay. What do you think?  
 306 *Speaker B:* Climb now, work later.

307 When it occurs as a logo on the aforementioned clothing, however, it cannot  
 308 be interpreted this way. There is not just one time which fixes the referent  
 309 of "now" and "later" in this case. Rather it expresses something along the  
 310 lines of "always go climbing before doing your work." In fact, it seems to have  
 311 much the same logical structure as a puzzling case noted (but not addressed)  
 312 by Kaplan:

313 (i) Never put off until tomorrow, what you can do today.

314 The difference, of course, is that (i) contains a lexical item "never" which  
 315 quantifies over temporal values, making the binding of "tomorrow" and  
 316 "today" explicit. But this difference is trivial—it is obvious enough that the  
 317 quantifier "always" is implicitly present as an unarticulated constituent of  
 318 (h). Thus we have a case of a **DC** in which the indexicals are forced to take  
 319 narrow scope with respect to an operator. This is a clear example of an **MO**.

#### 324 **Metalinguistic Monstrosity and Varieties of Quotation**

321 The above sections give us reason to accept that **DCs** are **MOs**. I will now argue  
 322 that, although they are **MOs**, they are not monsters. This leads naturally to the  
 323 conclusion that their context-shifting powers are not the result of any semantic

324 operation but is best thought of as the result of a pragmatic process. This will  
325 lead me to conclude that DCs pose no serious challenge to Kaplan's theory  
326 of indexicals and, in particular, to his restriction of the range of admissible  
327 contexts to proper contexts. Consequently, I shall argue that Kaplan is correct  
328 to classify sentences like "I am here now" as encoding contingent a priori  
329 truths.

330 That DCs are not monsters is fairly self-evident. Monsters are linguistic  
331 operators. No lexicalised operator can be discerned in the DCs discussed in  
332 the literature. It is not the addition of a new constituent to the sentence "I  
333 am here now" which performs the role of an MO when this sentence occurs  
334 within a DC. It is the particular use the sentence is put to. One could perhaps  
335 pursue the line that an unarticulated constituent is responsible for the MO  
336 and is thus a monster but there seems little evidence or motivation for such a  
337 view.<sup>6</sup>

338 If DCs are MOs but are not monsters, then they are most naturally under-  
339 stood as behaving exactly like the paradigmatic metalinguistic MOs recognized  
340 by Kaplan which were discussed earlier. Quotation, for example, is an MO  
341 because it can take an indexical sentence and block the indexicals within it  
342 from taking their customary wide scope positions. It does so by *mentioning*,  
343 rather than *using*, the sentence. I suggest that DCs are the result of MOs which  
344 do exactly the same thing. A DC is created by taking an indexical sentence and  
345 recording it ready to be mentioned in a new setting at a later date. Construed  
346 in this way, DCs are not constructed by a semantic operation. Thus the criti-  
347 cism levelled at Kaplan which draws on DCs as apparent counter-examples  
348 to his semantic theory for indexicals is misplaced. A distinctive feature of  
349 DCs is that they require a rich contextual setting. This contextual setting is  
350 not the minimalistic sequence of parameters required for saturating indexi-  
351 cal expressions, but a far wider notion of "context," incorporating various  
352 complex conventions surrounding human interaction. These conventions  
353 are essential to the performance of a DC. But such features are not semantic  
354 features. Thus it is natural to understand the MOs involved in generating DCs  
355 as pragmatically licensed, rather than semantic. I will now explain in detail  
356 the pragmatic process that I suggest is at work in these cases.

357 Quotation is surprisingly varied. Whereas it was once assumed that quota-  
358 tion is a simple device for self-nominalisation, enclosing a string of expressions  
359 within quotation marks to generate a name of that string, it is now widely

---

6 This point is argued for convincingly by Predelli (1996).

360 noted that quotation is not restricted to this simple operation. Consider a case  
 361 where it does seem to behave in such a way. For example:

362 (j) “Schnee” is a German expression which stands for snow.

363 In (j) the name “Schnee” names the expression “Schnee” but that expres-  
 364 sion is not used in any meaningful way in (j). We are simply exploiting the  
 365 convention whereby quotation marks name the expressions enclosed within  
 366 them. We could just as easily have exploited a different naming convention or  
 367 indeed stipulated one. For example, I can stipulate a convention for naming  
 368 an expression as follows:

369 (k) Let whichever German word stands for snow be called “Angelika.”  
 370 Angelika is often uttered by German speakers when in the presence of  
 371 snow.

372 The convention employed in (k) is perfectly clear and comprehensible. It does  
 373 not require any grasp of the German expression for snow or even an ability  
 374 to recognize that expression. We simply report facts about the expression  
 375 by utilising a (descriptively introduced) name of it. These sorts of examples  
 376 demonstrate that quotation behaves as a purely “mentioning” device in such  
 377 contexts. The semantic content of the expression that is being quoted is wholly  
 378 inert in these contexts (hence its unproblematic absence in (k)).

379 Other uses of quotation, however, are less simple. Newspaper headlines,  
 380 for example, commonly employ quotation not only to report speech but also  
 381 to convey information encoded by that speech. Here are a few examples taken  
 382 at random from the BBC News website on one visit:

383 Woman “killed dad and buried him”  
 384 Army ads “won’t appeal to new soldiers”  
 385 Financial services “pivotal to Brexit deal”  
 386 “Chronic” nurse shortage and Meghan “mania”

387 In all of the above cases, the quoted material is a speech report. However, there  
 388 is more going on here than *just* a speech report. Compare it to the following,  
 389 more straightforward speech report, taken from Rap Artist Chuck D’s (1998,  
 390 193) autobiography:

391 Ice-T was in the video because I saw him while he was in Arizona  
 392 and asked him if he wanted to be in the video. He said, ‘Cool.’



393 This is simply a speech report—reporting the words used to accept an offer  
394 to feature in a promotional video. In the previous examples, however, the  
395 direct quotation is not simply a report of the words used by whoever uttered  
396 them—it also draws attention to the speech act they were used for and, in  
397 doing so, draws attention to (and *uses them to express*) their content. In this  
398 usage, which is very common in newspaper headlines, the quotation does not  
399 just name the expressions used but also establishes that they were used to  
400 *allege* something. In such cases where an allegation is reported, the reporter  
401 does not implicate herself as one making the allegation, she simply reports  
402 the allegation and reports that such an allegation has been made. But we can  
403 also find cases where quotation is employed not just to report a claim, but  
404 also to *endorse* that claim:

405 (1) Kaplan’s example of a kidnapped heiress, locked in the trunk of a car,  
406 who has lost all track of time and of her location, yet who can still  
407 think to herself “it is quiet here now,” demonstrates clearly that, “igno-  
408 rance of the referent does not defeat the directly referential character of  
409 indexicals.”

410 In (1) the quoted material (from [Kaplan 1989a, 536](#)) at the end of the passage  
411 both reports Kaplan’s view and, at the same time, endorses it. The quoted  
412 material is not simply named; its content is asserted. Quotation of this sort,  
413 labelled “mixed quotation” by Cappelen and LePore (1997) because of its dual  
414 role as reported and asserted, is quite common.<sup>7</sup> Récanati (2010) helpfully  
415 distinguishes between “closed” and “open” forms of quotation to make sense  
416 of the distinctions at play in these cases. Following Davidson (1979), Récanati  
417 interprets quotation marks as performing a demonstrative role whereby the  
418 linguistic material (L) quoted is demonstrated as the referent of the quoted  
419 expression “L” in reports of the form “S said that ‘L.’” The difference between  
420 open and closed quotation is that closed quotation recruits the demonstrated  
421 material to play the syntactic role of a singular term, whereas open quotation  
422 is any form of quotation that cannot be so construed:

423 The contrast between open and closed quotation is illustrated by  
424 the following pair of sentences:

---

7 Davidson (1979, 29) first drew attention to such “mixed case[s] of use and mention.”

- 425 (7) Stop that John! ‘Nobody likes me’, ‘I am miserable’ ... Don’t  
 426 you think you exaggerate a bit?  
 427 (8) John keeps crying and saying ‘Nobody likes me’.

428 In (7) a token of ‘Nobody likes me’ and ‘I am miserable’ is dis-  
 429 played for demonstrative purposes, but is not used as a singular  
 430 term, in contrast to what happens in (8), where the quotation  
 431 serves as a singular term to complete the sentence ‘John keeps  
 432 crying and saying \_\_\_\_’. Sentence (7), therefore, is an instance  
 433 of open quotation, while (8) is an instance of closed quotation.  
 434 (Récanati 2010, 231)

435 If we acknowledge this distinction, we ought to agree with Récanati that  
 436 there is a fundamental difference in linguistic role between the demonstrated  
 437 linguistic material in closed quotation and that in open quotation. Whereas  
 438 closed quotation recruits (a token of) the linguistic material as a singular  
 439 term which is naturally understood as referring to itself (as a type in most  
 440 instances), open quotation does not feature any singular term which naturally  
 441 presents itself as requiring a referential interpretation. Récanati’s proposal is  
 442 that the sense in which the quoted material acts as a *demonstration* in open  
 443 quotation is wider than the customary sense in which demonstratives are  
 444 taken to have their reference fixed by an accompanying demonstration. In  
 445 open quotation, quoted material *demonstrates* in the sense of providing a  
 446 performance or picture that represents through a form of mimicry. Thus, in  
 447 Récanati’s example (7) above, John’s speech is quoted as a means of mimicking  
 448 his self-indulgent utterances.

449 Understood as mimicry, open quotation has monstrous potential. This is  
 450 unsurprising, of course: we have already noted that quotation is a metalinguistic  
 451 monster. Mimicry is clearly an attempt to represent a content which is,  
 452 in some sense, displayed from a perspective other than that of the speaker,  
 453 namely the perspective of the one that the speaker mimics. Consider my  
 454 report of my teenage daughter’s recent request for a new pair of trainers:

- 455 (m) Amy has asked me to get her a “sick” new pair of trainers.

456 I do not, ordinarily, use the word “sick” with the sense it has been assigned  
 457 in (m). If I am honest, I confess that I am not entirely sure what the exten-  
 458 sion of the term “sick” is when my daughter and her friends use it in their  
 459 conversations. But I know enough about it’s meaning to know that it is a

460 desirable property of footwear in my daughter's opinion (hence this use of  
461 the expression is not synonymous with other, more common, uses of "sick"  
462 in English) and this is readily communicated (to others who have at least the  
463 same level of acquaintance with this term as I do) by my utterance of (m).  
464 But I do not communicate *my* judgement as to the possession of this property  
465 by any footwear when I use this term in (m); I communicate my daughter's  
466 judgement. Thus we have a form of context-shifting operator present in (m).  
467 The quoted expression "sick" shifts the context to one in which Amy judges  
468 things to have the property that she takes that expression to encode. As Ré-  
469 canati (2010, 260) notes this may not amount to a full-blown MO as it is not  
470 clear that shifting from my idiolect to Amy's is best represented by a shift in  
471 the sequence of parameters we would normally associate with a linguistic  
472 (as opposed to metalinguistic) context. Nonetheless, it illustrates nicely the  
473 potential that open quotation has for shifting the perspective away from that  
474 of the speaker in a way that is commonplace in ordinary discourse. Indeed we  
475 can extend the usual notion of a linguistic context to accommodate such met-  
476 alinguistic operations easily enough, by incorporating an "idiolect" parameter  
477 for the context (see Récanati 2010, 260), according to which my utterance of  
478 (m) will be interpreted as employing a context-shifting operator (quotation  
479 marks) to recruit the content of the expression "sick" as assigned in Amy's  
480 idiolect to act as a constituent of a proposition whose content is otherwise  
481 assigned in accordance with my idiolect.

482 Open quotation, understood thus, has a number of advantages, most no-  
483 table among them being the fact that we can simultaneously maintain that  
484 quotation is (i) metalinguistic, (ii) an MO, and (iii) *used* (rather than merely  
485 mentioned) as a means of communicating information. Furthermore, as I will  
486 now illustrate, it provides a perfect explanatory model for the MOs discussed  
487 in this paper. My proposal is that DCs are best analysed as akin to instances  
488 of open quotation. It follows that DCs are MOs but this poses no threat to  
489 the Kaplanian claims that there are no linguistic monsters in English, only  
490 metalinguistic MOs, and that DCs do not provide counter-examples to the  
491 proposal that there are special logical truths of indexicals languages or to any  
492 of Kaplan's proposed logical truths of indexical languages.

493 To see how MOs can be interpreted on the same lines as instances of open  
494 quotation, it will be helpful to consider a range of similar phenomena involv-  
495 ing intuitive context-shifting operations.

496 *Expressives* are expressions which encode a speaker-attitude alongside, but  
 497 independently of, their truth-conditional content (if they have any).<sup>8</sup> If I utter  
 498 the expression “yummy,” upon encountering a delicious foodstuff, I express  
 499 my positive attitude towards its flavour. But now consider the phenomenon  
 500 of child-directed speech. In the years before my daughter was old enough to  
 501 desire “sick” new trainers, I used to feed her baby food: pureed vegetables,  
 502 rusks, and other assorted foodstuffs which I personally do not find even  
 503 remotely appetising. Yet, it was common for me to feed her and to accompany  
 504 the process with enthusiastic utterances of “yummy.” Intuitively, I was not  
 505 expressing my positive attitude towards the taste of the food; rather I was  
 506 expressing (or perhaps encouraging) my daughter’s positive attitude. A simple  
 507 explanation of what is happening here is that the expressive “yummy” encodes  
 508 the attitude of the speaker. But in cases like this, we have an implicit open  
 509 quotation operator, which effects a context shift from speaker-attitude to the  
 510 attitude of the quoted speaker. This is realised, as in the cases considered above,  
 511 by an act of mimicry. By mimicking the reaction (or, perhaps, *desired* reaction  
 512 in this case) of my daughter, I implicitly apply a form of open quotation to  
 513 her utterances.

514 We encounter the same thing when we consider expressives with a truth-  
 515 conditional component. The expression “eurocrat” is a mild, and slightly  
 516 comical,<sup>9</sup> pejorative expression used by anti-EU British politicians (and those  
 517 who support them in this regard) to denote the politicians and their fellow  
 518 officials who form the European Parliament and administer the bureaucracy of  
 519 the European Union. An utterance of the sentence, “I think it is hilarious that  
 520 Farage has to spend his time hanging out with all those eurocrats,” made by  
 521 someone who obviously does not share Farage’s attitude of contempt towards  
 522 the bureaucrats in question, is naturally read as an open quotation which  
 523 would most perspicuously be represented as such:

524 (n) I think it is hilarious that Farage has to spend his time hanging out with  
 525 all those “eurocrats.”

---

8 “Pure” expressives like “ouch” and “oops” appear to make no truth-conditional contribution to utterances and to simply encode a speaker attitude, whereas, e.g., pejoratives (including racial or sexual pejoratives) are often thought to encode both a truth-conditional content (an extension, namely those who the speaker intends to denote by the term) and a speaker attitude (of derogation towards the members of that extension).

9 Admittedly, its comedic quality has been somewhat diminished since the UK’s recent referendum result.

526 Interpreted thus, it is obvious that the pejorative force of the expression “eu-  
527 rocrat” has shifted away from the speaker-attitude to the attitude of Farage.  
528 The use of the term is again interpreted as mimicry.

529 *Predicates of personal taste* have also been noted as displaying similar be-  
530 haviour. If Mary says “Rollercoasters are fun” and John says “Rollercoasters  
531 are not fun,” Mary and John are disagreeing *faultlessly*. That is to say that  
532 although it appears that one is asserting a proposition while one asserts the  
533 negation of that proposition, there is a sense in which both are speaking the  
534 truth, without either speaker misrepresenting the facts. Semantic Relativists  
535 like Lasersohn (2017) explain this by taking the truth of their utterances to be  
536 sensitive to a contextual parameter included in the circumstance of evalua-  
537 tion which ensures that the standard for truthful attribution of fun may differ  
538 between the two utterances. But now consider the case where Mary asks John,  
539 immediately following his rollercoaster ride, “Was that fun?”. Intuitively, the  
540 relevant standard here is not Mary’s but John’s. She is asking if *he* found it fun.  
541 Again, we can make sense of this scenario by understanding the evaluation of  
542 the attribution of the property encoded by “fun” to be relative to a parameter  
543 which is usually set to the speaker of the expression but in cases like this is  
544 shifted to the addressee. Again this can be understood as resulting from an  
545 implicit open-quotation device to mimic the addressee of the question. Note  
546 that mimicry here does not have to be a convincing performance, it simply  
547 needs to present the attitude or perspective of the target agent to whom the  
548 attitude is being attributed. I can do the same thing when I feed my pet guinea  
549 pigs some dried pellets of food and ask them, “Is that tasty?”. I do not need to  
550 be providing a convincing impression of a guinea pig to make it clear that the  
551 relevant standard of, and perspective on, tastiness here is that of my guinea  
552 pigs (or, at least, that which I attribute to them), not mine.

553 The above examples demonstrate that context-shifting is familiar for a range  
554 of expressions.<sup>10</sup> What then of the content of indexicals? Can we provide

---

10 Of course, not all will share my analysis of these cases as instances of open quotation. For example, irony of the sort displayed in examples like (m) and (n) may inspire competing analyses from Griceans. I do not have space here to mount a detailed defence of my analysis of irony and related phenomena, but hope to have made it clear that the analysis is a plausible one for a range of phenomena that are importantly similar to the cases we are concerned with. As well as drawing on Récanati’s approach, my analysis has some similarities with the echoic analysis of irony and related phenomena adopted in Relevance Theory (see, e.g., Wilson 2006). The Relevance Theoretic approach is applied by Bianchi (2014) to echoic uses of slurs. An important point, emphasized by her, is that when we echo or imitate the perspectives of others, we do not have to extend the echoing to the whole content of an asserted proposition—we often only express a

555 examples where the same context-shifting operation shifts the reference of  
 556 an indexical in the way that we expect MOs to do? In fact such examples are  
 557 easy enough to find. First consider another example of child-directed speech.  
 558 A nursery teacher, talking to a very young child who had her birthday the  
 559 day before says: “Did mummy and daddy take you somewhere nice for your  
 560 birthday?” Not only is the contextual standard for “nice” shifted to that of  
 561 the addressee, but also the content of the terms “mummy” and “daddy” have  
 562 shifted. These expressions behave very much like indexicals in that when  
 563 uttered without qualification, they refer to the parents of the speaker. But here,  
 564 the only qualification arises as a consequence of the nature of the context.  
 565 That context generates a construction best understood as an open quotation:  
 566 “Did ‘mummy’ and ‘daddy’ take you somewhere ‘nice’ for your birthday?”  
 567 where the open quotation operation shifts the context away from that of the  
 568 speaker to the addressee in order to fix the content of the quoted expressions.

569 Of course, it might be replied that this example can equally be explained  
 570 by appeal to ellipsis. It might be thought that the indexical-like features of  
 571 “mummy” and “daddy” are best explained by appeal to an elided possessive  
 572 marker “ $\alpha$ ’s mummy,” which may be an obviously indexical possessive pro-  
 573 noun “my mummy,” “your mummy,” “her mummy,” etc. Be that as it may,  
 574 there are other cases which make it perfectly clear that indexicals can be  
 575 shifted by open quotation. Indeed we saw one above from Récanati, which I  
 576 repeat here:

577 (o) Stop that John! “Nobody likes me,” “I am miserable” ... Don’t you think  
 578 you exaggerate a bit?

579 It is clear that the referent of the “me” and “I” in (o) is not the speaker of  
 580 (o) but the person they are mimicking, namely John. Examples like (o) are  
 581 not uncommon and are a clear example of the use of open quotation as an  
 582 MO that shifts the context that the indexicals contained within it are indexed

---

perspective distinct from our own with regard to a *part* of that proposition. For example, Bianchi interprets the sentence “As I reached the bank at closing time, the bank clerk helpfully shut the door in my face” as containing an instance of echo or imitation only with regard to the expression “helpfully” (2014, 39). This is the same feature that I am appealing to open quotation to explain in many of the examples above. Bianchi draws on this analysis to explain seemingly non-offensive uses of slurs, such as we see in appropriation (cases where the usual targets of a slur use the expression in a way that removes its derogatory aspect). While I agree that there are echoic uses of slurs such as (n), I would not extend this analysis to appropriation (my own account of appropriation can be found in Scott and Stevens 2019); other examples of echoic uses of slurs and of expressives more generally are discussed in Stevens and Duckett (2019).

583 away from the parametric settings of the overall sentence to another context  
 584 for those quoted segments of the utterance. Using subscripts to display the  
 585 relevant indexes, the situation is something like this:

586 [Stop that John! “[Nobody likes me]<sub>c<sub>2</sub></sub>,” “[I am miserable]<sub>c<sub>2</sub></sub>” ... Don’t  
 587 you think you exaggerate a bit?]<sub>c<sub>1</sub></sub>

588 John is the addressee of  $c_1$ , and not the agent; he is the agent of  $c_2$ , not the  
 589 addressee.

590 Notice that open quotation is being considered as an explicit (albeit metalinguistic)  
 591 operation in the above examples. The operator (quotation marks)  
 592 is ambiguous between open and closed quotation producing functions but  
 593 it is explicit in the syntax of the written language. But, of course, utterances  
 594 are not always inscribed. Except in rare cases where quotation marks are  
 595 “signed” by a gesture which conventionally signals that the words uttered  
 596 contemporaneously with that gesture are being quoted, it is up to hearers to  
 597 identify quotation from features of the context. Pragmatic aspects of utterance  
 598 interpretation come to the fore in such situations. Consider the following  
 599 pairs of utterances:

- 600 (p) That guitarist, John, is performing tonight.  
 601 (q) That guitarist, John, who can’t actually play the guitar to save his life, is  
 602 performing tonight.\*

603 The awkwardness of (q) (indicated by the \*) follows from the apparent contradiction  
 604 which results from simultaneously describing John as bearing a  
 605 property and then denying that he bears that very property. Were one to hear  
 606 an utterance of (rather than read an inscription of) (q), however, one would  
 607 most likely apply a principle of charity and interpret the utterance in a way  
 608 which resolved this potential infelicity, such as (q<sub>1</sub>):

- 609 (q<sub>1</sub>) That “guitarist,” John, who can’t actually play the guitar to save his life,  
 610 is performing tonight.

611 In other words, the term “guitarist” is interpreted as being subject to an open-  
 612 quotation operator, shifting its usual extension to one that includes John (who  
 613 is exempt from the extension of the standard English term). It may be read as  
 614 synonymous with “so-called ‘guitarist,’” hence behaving much like the term  
 615 “sick” discussed above: the target of a metalinguistic operator that shifts the  
 616 idiolect (or other metalinguistic feature) according to which it is interpreted.

617 This process of pragmatically guided utterance interpretation need not  
618 apply only to subsentential elements but can equally be applied to whole  
619 sentences. A few years ago, there was something of a craze for purchasing  
620 audio recordings to be played as an answerphone message. A popular recording,  
621 used as an example above, was the snippet of Arnold Schwarzenegger's  
622 character from the *Terminator* movies uttering the line "I'll be back." If I  
623 call my friend and hear this message, I do not interpret it as expressing the  
624 proposition it was originally used to encode. I interpret it as saying that the  
625 person who I have called is temporarily absent and soon to return. It is, in  
626 fact, interpreted as if the person being called were able to respond to my call  
627 from their current location and say "As Arnold Schwarzenegger says: 'I'll be  
628 back!'" In other words, I understand the utterance as displaying the utterance  
629 made by Schwarzenegger and recruiting it to communicate information. I  
630 interpret it in precisely the same way as an instance of open quotation. The  
631 message is interpreted in just the same way as I would interpret the utterance  
632 of my friend who explicitly mimicked Schwarzenegger's character, monotone  
633 pronunciation (perhaps even accompanied by distinctive bodily movements)  
634 and all, when in my presence.

635 Of course, mimicking an iconic actor or fictional character, by uttering  
636 an iconic line from an iconic movie is one thing, but what about ordinary  
637 answerphone messages, written notes, etc.? Who, or what, is being mimicked  
638 in these cases? Mimicry in these cases is more mundane but best understood  
639 as mimicry nonetheless. All that happens in these cases is that the speaker  
640 mimics *themselves* saying what they would say, were they able to inhabit the  
641 impossible (that is, improper) contexts they would need to be in to otherwise  
642 communicate this information. That is to say, when one needs to communicate  
643 information from a context unavailable to them, one must find an alternative  
644 method of relaying the information. By preparing in advance a message to be  
645 retrieved by ones intended audience in this context one is able to overcome this  
646 obstacle. But this is achieved not, as is often assumed, by somehow making  
647 an utterance "from a distance" but by recording in one form or another an  
648 instance of oneself performing the speech act one would want to make at  
649 that context if able to, ready to be displayed there. In doing so, one does  
650 not encode the proposition that would be obtained by pairing the uttered  
651 sentence with the improper context in question, but simply prepares a string  
652 of linguistic material that mimics the intended performance and then exploits  
653 the various media which permit this mimicry to be planted in advance ready  
654 to be deciphered when encountered. What is deciphered is not an utterance



655 by an absent agent, it is a previous utterance deliberately placed in a situation  
656 where it will pragmatically trigger a process of interpretation precisely akin  
657 to that by which we read utterances of open quotation as discussed above.  
658 There is no fundamental difference between my employing open quotation to  
659 use an instance of the indexical “I” to pretend to be someone I am not, and  
660 my employing open quotation to use an instance of the indexical “here” to  
661 pretend to be speaking at a location where I am not. DCs are not utterances  
662 made at a distance, they are recorded performances pragmatically recruited  
663 to mimic intended and otherwise impossible utterances. This distinction is  
664 far from trivial: it demonstrates that DCs are generated by pragmatic features  
665 of communication and are thus not data to be accommodated by semantic  
666 theory.

667 The only difference between the uses of open quotation that I have discussed  
668 above and the full-blown DCs is that the syntactic role of quotation to mark off  
669 the shift will clearly be uncalled for in the latter case. Accordingly, we should  
670 not expect quotation to be readily recoverable in a case where a DC consists  
671 entirely of a mimicked performance, whereas it is essential for embedded  
672 occurrences like we see in “Amy has asked me to get her a ‘sick’ new pair of  
673 trainers.” Only when DCs are embedded would explicit quotation marks be  
674 felicitous, as we saw with “As Arnold Schwarzenegger would say: ‘I’ll be back.’”  
675 Other DCs would need to be placed in similarly embedded constructions to  
676 achieve the same result. For example: “If I could speak at the context where  
677 you will hear this message I would report that, ‘I am not here now,’” and so  
678 on.

## 675 **Semantics, Pragmatics, and Displaced Communications**

680 In this section I want to briefly say a few things in defence of the view presented  
681 above and to point out its advantages over competing accounts of DCs. In  
682 recent years the most vocal and influential critic of Kaplan’s account of logical  
683 truth for indexical languages has been Stefano Predelli. Although couched  
684 within a position sympathetic to Kaplan’s semantic project, Predelli takes issue  
685 with Kaplan’s decision to limit the possible combinations of sentences with  
686 contexts to proper contexts. Drawing on the DCs discussed above, Predelli  
687 argues that Kaplan is simply wrong to assume that utterances require their  
688 agents to be located at the times and places they occur. Furthermore, he offers  
689 an ingenious proposal as to how an extension of Kaplan’s theory to include  
690 improper contexts can be motivated and put into practice.

691 According to Predelli (most comprehensively in 2005) no semantic theory  
 692 for indexical languages can be complete unless it has the resources to ac-  
 693 commodate the role of speaker intentions in fixing the parametric settings of  
 694 indexicals as uttered. In particular, Predelli maintains that speaker intentions  
 695 are crucial to determining the context with which an indexical sentence must  
 696 be paired in order to correctly model the actual utterance. Presented thus,  
 697 Predelli's position may not sound particularly distinctive—challenges to the  
 698 attempt to model meaning by formal semantics alone without recognition of  
 699 the role played by the speaker intentions behind the utterances whose mean-  
 700 ing we are attempting to model are common from those who maintain that a  
 701 theory of pragmatics is needed to explain linguistic meaning. The novelty of  
 702 Predelli's position however rests on his desire to reconcile his approach with  
 703 a philosophy of language that assigns the core role of explaining meaning  
 704 to formal semantics. To bring about this reconciliation of speaker intentions  
 705 and semantically assigned meanings, Predelli draws a distinction between the  
 706 workings of a formal system which calculates truth-conditions for utterances  
 707 and a “pre-semantic” arena in which the inputs to this system must be first  
 708 determined. It is in this latter area that speaker intentions become significant.  
 709 Before I can employ a formal system to calculate the truth-conditions of a  
 710 speaker's utterance, I first must determine *which* utterance she has made  
 711 (which *proposition*, in other words, she has said). To take a trivial example not  
 712 involving indexicality, I can only know the truth conditions of an utterance  
 713 of the sentence “John is sitting beside the bank” if I know which lexeme the  
 714 ambiguous English word “bank” encodes in that sentence. So I must deter-  
 715 mine which lexeme the speaker intended before I input her utterance into  
 716 the formal semantic theory which then returns its truth-conditions as output.  
 717 The same thing happens when one utters an indexical sentence, according to  
 718 Predelli, but now the pre-semantic task is to determine the correct context  
 719 that this sentence must be paired with and, as with the case of disambiguating  
 720 a lexical ambiguity, the only correct answer here will be that which identifies  
 721 the speaker's intention.

722 Certainly there is much to agree with in Predelli's account, and its sub-  
 723 tleties are not always recognised by his critics. For example, the objection often  
 724 directed at Predelli that he is guilty of Humpty-Dumptyism<sup>11</sup> is misguided.  
 725 Humpty-Dumpty does not resolve ambiguities or select contexts to pair in-

---

11 “Humpty-Dumptyism” is the pejorative term for a semantic theory of the bizarre and implausible sort envisaged by Lewis Carroll's fictional character who insisted that his words mean simply whatever he wants them to mean. Responses to the charge are given in Predelli (2011).

726 dexicals with, he simply rejects on whim the existing semantic assignments  
727 given to the elements of his vocabulary and selects alternative ones again  
728 on whim. Humpty's speech is thus effectively unreadable to any semantic  
729 system employed by anyone other than himself. He scrambles the inputs to  
730 semantic theory into a code known only to himself. Thus communication  
731 breaks down. Applying this analogy to indexical terms, Humpty would be  
732 guilty of modifying the characters of indexicals at whim. We would not know  
733 which function from contexts to contents was encoded by his use of, e.g., "I"  
734 and hence could not calculate its content. But, as Predelli (2005, 58) explicitly  
735 states, characters are left untouched by the pre-semantic task of sentence-  
736 context pairing. Agents of answerphone messages are not using "I" to pick  
737 out anyone other than the agent of the context, they are simply selecting an  
738 improper context to pair their use of the term with.<sup>12</sup>

739 Predelli will therefore reject my interpretation of DCs as MOs.<sup>13</sup> There is  
740 no operator on character according to his view, only a selection of a context  
741 that we have been wrongly denied in Kaplan's theory. Dropping Kaplan's  
742 restriction to proper contexts allows us to accommodate DCs as respectable  
743 utterances made "at a distance." What, then, is to be said in favour of my  
744 view over Predelli's? I think that Predelli's view, for all its ingenuity, suffers a  
745 number of drawbacks that my account is not prone to.

746 Firstly, as argued in detail by Stevens (2009), Predelli's position stands or  
747 falls on the strength of the intuition that DCs really are cases where utterances

---

12 An alternative source of the humpty-dumpty objection to Predelli that I have heard attacks the account on the grounds that it allows the speaker to pair a sentence with any context that they choose, hence their choices as to that pairing could, in principle, be just as private to them as Humpty's choices about meaning assignments are to his idiolect. I don't find this objection compelling—there is no obstacle to Predelli admitting that there are success conditions placed on successful communication that apply to the pairing of sentence and context just as there are for resolving lexical ambiguities. I can successfully encode a number of things by "bank," but not just *anything*. My intended meaning will only succeed if it conforms to existing conventions about English usage. Similarly for the intentions I have about the contexts I pair my sentences with.

13 Predelli (1996) addresses the relation between monsters and DCs to some extent. Although in this paper Predelli does not consider all alleged cases of monsters, focusing solely on the famous "never put off until tomorrow what you can do today," he explicitly appeals to his intentionalist framework to explain away the apparent monstrosity of this example by maintaining that it should be understood as encoding multiple intended DCs. Discussion of his account of this particular sentence takes us beyond the scope of this paper. For other interesting discussions of cases that seem to involve one and the same sentence expressing multiple DCs as it is decoded repeatedly, see Egan (2009) and O'Madagain (2014). See Predelli (2014) for further discussion of monsters.

748 are made at a distance.<sup>14</sup> But this intuition is fragile and sensitive to varying  
 749 examples. The presentation of the data makes a difference to the intuition.  
 750 For example, while it is true that when I phone Amy and hear her recorded  
 751 answerphone message “I am not here now, please leave a message,” I understand  
 752 that this expressed a prior intention on her part to communicate to  
 753 anyone who hears the message the fact that she is not present at that later  
 754 time, there are features one would expect an *utterance* to have which are  
 755 lacking in this scenario. For example, it would be very odd of me to accuse her  
 756 of lying, or even unwittingly telling an untruth, if I knew that she was in fact  
 757 present at the location of her answerphone when I called her. The appropriate  
 758 things to say in such a case would be something like “your *answerphone* is  
 759 wrong/misleading/in need of updating, etc.” not “you are wrong...”.<sup>15</sup> Similar  
 760 points can be made about all DCs. We intuitively recognize a gap between  
 761 these devices of communication and ordinary utterances, but this gap goes  
 762 intrinsically unrecognised on Predelli’s intentionalist account, according to  
 763 which the DC is a straightforward utterance.

764 Secondly, Stevens (2009) also points out that the intuition that a DC is a  
 765 genuine utterance appears hard to reconcile with the equally strong intuition  
 766 that an utterance is made at the time of encoding of the message. To, as Sidelle  
 767 (1991, 535) puts it, engage in “arranging to make an utterance at a later time,  
 768 or, if one likes, deferring an utterance,” is not to engage in uttering something  
 769 while making those arrangements. This is especially clear in Predelli’s case.

---

14 Predelli (2011) responds to the objections raised by Stevens (2009) by characterizing those objections as founded on a mistaken conception of the proper role of semantic theory. On Predelli’s characterization, Stevens is denying that “the evidence put forth by true instances of ‘I am not here now’ [should] constrain the shape of an empirically adequate semantic account” (2009, 301). I agree that this would be a mistaken view of the role of semantic theory; however, it is a misrepresentation of the objection from Stevens, who clearly rejects the intuition that there *are* any true instances of “I am not here now” to be accounted for. Of course, Predelli is correct to note that the question of where the line should be drawn between semantic theory and pragmatic theory is a controversial one. I take the considerations in [this section](#) and the [preceding one](#) to lend compelling support to views like that forwarded in Stevens (2009) and Récanati (2010), according to which that line is decisively drawn in a way that makes pragmatic theory responsible for explaining DCs rather than semantic theory as Predelli maintains.

15 As a referee pointed out to me, intuition seems to shift back in the other direction in a case where there is a deliberate deception. Suppose that Amy does not want to speak to me and deliberately leaves her answerphone on so that I will think she is not there. Now the intuition that she is lying has more traction. I take this point. However, I am content to use the example to illustrate that our intuitions are unstable—the intuition that agents make utterances at places and times other than where and when they are situated is malleable in a way that the intuition that agents make utterances in Kaplanian proper contexts is not.

770 If the context of utterance is the intended context of utterance, then I utter  
771 nothing at all when recording my answerphone message; I simply get things  
772 ready for an utterance to occur later on. Perhaps this is so, but insofar as the  
773 position is motivated by our intuitions regarding DCs, this counterintuitive  
774 consequence counts against Predelli's account. On my account, however, it  
775 can be explained easily enough. One is simply engaging in an act of pretense  
776 when recording the message, mimicking what one would say if located at the  
777 time and place of the context at which our intended audience will hear our  
778 performance. Unlike Predelli's account, this entails no claim about utterances  
779 being displaced from their proper contexts. The only utterance that takes place  
780 is that which is made when recording the message, although it is not uttered  
781 with assertoric force; it is simply the product of an act of mimicry, ready to  
782 be displayed in a different context. The underlying intuition that motivates  
783 Predelli's intentionalist account, namely the intuition that I am deliberately  
784 aiming to communicate things at contexts other than the one in which I am  
785 located, is preserved without endorsing the counterintuitive consequences of  
786 construing this as a form of utterance at a distance, by accommodating that  
787 intuition within a purely pragmatic explanation.

788 Technology can open up the possibility of previously outlandish uses of  
789 language but this is best explained pragmatically, not through a reconstruction  
790 of an otherwise perfectly acceptable semantic theory. Consider a recent  
791 technological advance which facilitates an unusual application of indexicals:  
792 The Rock Group, *Dio*, recently performed a series of concerts in which their  
793 deceased vocalist Ronnie James Dio was replaced on stage by a hologram.  
794 The hologram appears to be singing as it mouths along to pre-recorded vocal  
795 tracks from Ronnie. This holographic rendition of Ronnie, convincing though  
796 it may appear, is not of course really singing. The hologram is not causing  
797 any vibrations in the air, picked up by a microphone, etc. It is just a visual  
798 representation of a dead person, carefully synchronised with recordings of  
799 that person's voice. But we can exploit this pretense to the full. When per-  
800 forming in London, we can make our holographic Ronnie "say" things (i.e.,  
801 mouth along to recorded utterances of Real Ronnie's) like "it is great to be in  
802 London tonight!". Of course, Ronnie himself is not saying anything; Ronnie,  
803 unfortunately, is dead. Suppose that on the evening of 20<sup>th</sup> December 2019,  
804 the hologram is made to "say" this sentence: "It is cold tonight in London!".  
805 Has the proposition that London is cold on the evening of 20<sup>th</sup> December 2019  
806 been expressed? It seems reasonable to agree that it has, although it is equally  
807 obvious that Ronnie was not the agent who expressed that proposition (the

808 most likely agent, or agents, would be those responsible for generating and  
 809 controlling the hologram). But what about a case where Holographic Ronnie  
 810 “says” something using the first-person pronoun like “I am so happy to be  
 811 here with you in London tonight!”. Again, I think it is obvious enough that  
 812 Holographic Ronnie has not said anything (not being an agent, he cannot be  
 813 the agent of an utterance after all). Nor, for that matter, has Real Ronnie said  
 814 that he is happy to be in London on the evening of 20<sup>th</sup> December 2019 (not  
 815 being alive, he is not able to be an agent and hence not able to be an agent of  
 816 an utterance). But this seems hard to square with Predelli’s view, according  
 817 to which the parametric settings that determine the content of an utterance  
 818 are fixed by intentions. In this case there clearly is an intention, just not an  
 819 intention on the part of Ronnie (Holographic or Real). But whoever produced  
 820 the hologram (let us assume it is a single individual for simplicity’s sake) had  
 821 an intention to combine a recorded utterance of Real Ronnie’s (perhaps recon-  
 822 stituted from several samples taken from previous utterances and hence not  
 823 identical with any one previous actual utterance) with a visual representation  
 824 of Real Ronnie to ensure that Holographic Ronnie “said” that he was happy  
 825 to be in London on the evening of 20<sup>th</sup> December 2019.

826 There does not seem to be any semantic difference between what is happen-  
 827 ing here and what is happening if I attached a note authored by Jones which  
 828 says “I am not here” to Smith’s door to express Smith’s absence. The agent of  
 829 the note, on Predelli’s account, is presumably Smith because it is his absence  
 830 I intend to communicate. Accordingly, Predelli’s account predicts that the  
 831 agent of Holographic Ronnie’s “utterance” is either Holographic Ronnie or,  
 832 perhaps, Real Ronnie (depending on which of these two, if distinguished,  
 833 the producer of the holographic performance intended). This, I think, can-  
 834 not be the right thing to say in this situation. No amount of intention can  
 835 make dead people agents of utterances after their death.<sup>16</sup> Surely what we

---

16 A similar objection to Predelli is raised by Sherman (2015, 594) who notes that Predelli’s intentionalist account makes apparently correct predictions about the cases where we have some choice over our use of indexicals, but struggles to explain cases where we don’t. The comment is made in passing but I assume he has in mind cases like this: a recently released addition to a range of ice cream has “I am vegan!” written on it. Predelli seems to have a simple explanation of what the “I” means here—whoever wrote this intended that it be paired with a context in which the ice cream is the agent. But now consider a case where I stand next to the freezer in the supermarket shouting, “I am vegan!” to passing shoppers, while intending the sentence I emit to be paired with a context in which the ice cream is the agent. Clearly, my intention will not be fulfilled. This suggests that there is more to the successful case than just the intentions of whoever produced the communication.

836 have is simply a case where someone is doing an extremely sophisticated job  
 837 of *pretending* that Ronnie James Dio is present and performing on stage by  
 838 displaying recordings of his previous speech in an act of mimicry. This, in my  
 839 view, is what we find in all cases of DCs. Ordinary utterances with Kaplanian  
 840 meanings are employed to allow us to pretend to say (or pretend that others  
 841 are saying) things unavailable to us when the utterances are construed literally.  
 842 There is no need to modify our semantic theory to accommodate a theory of  
 843 pretence.<sup>17</sup>

844 One final approach to the answerphone problem that seeks to accommodate  
 845 DCs within a wholly semantic framework is suggested (though not endorsed)  
 846 by Parsons (2011). The view merits brief consideration here as, again, it shares  
 847 some similarities with my proposal but the differences are significant. I have  
 848 argued that apparently true instances of “I am not here now” etc., are not  
 849 really true. They are false utterances made by speakers who utilise a pragmatic  
 850 process to facilitate their non-literal interpretation as *pretences*. Speakers are  
 851 relying on context to allow them to mimic utterances at different temporal or  
 852 spatial locations (or even by different agents) because that context will make  
 853 those shifted contexts salient (most routinely because those shifted contexts  
 854 are the ones that the hearer will be in when they decode the utterance).

855 It is helpful to compare my view to a radical form of what we might call “con-  
 856 tent relativism” (CR). CR is the position whereby the content of an utterance  
 857 is subject to modification depending on the context in which the utterance is  
 858 assessed (rather than the context in which it is *uttered*).<sup>18</sup> In Kaplanian terms,

---

17 The idea that DCs can be explained as pretend utterances is also defended by Voltolini (2006) and Connolly (2017). Voltolini’s strategy is to situate his explanation within a fictionalist semantics, while Connolly shares my preference for a pragmatic approach. I have a great deal of sympathy with Connolly’s approach which construes DCs as produced by participants knowingly and deliberately entering into a game of “externally-oriented make believe” (2017, 616). However, while I think our approaches are in the same vein, I think the situation he describes must be supplemented by the sort of analysis I propose if it is to explain the monstrous quality of DCs. For example, I have argued that binding of indexicals (Climb *now*, work *later*), and embedded context-shifting (Amy has asked me to get her a “sick” new pair of trainers) have important similarities with DCs that require the sort of approach I am urging.

18 CR is a more radical theory even than the controversial forms of semantic relativism (or, as it is sometimes called “truth relativism”) developed most notably by Lasersohn (2017) and MacFarlane (2014). Semantic relativism holds that truth is sensitive to context of assessment; CR holds that *what is meant* by an utterance depends on the context in which that utterance is assessed. Despite the clear logical space for CR to exist within any semantic framework which admits both contexts of utterance and contexts of assessment, few have been persuaded that CR is worth exploring. A rare exception (in addition to Parsons 2011, discussed shortly) is Weatherson (2009).



859 the context in which an utterance is assessed is the circumstance of evaluation.  
860 Whereas Kaplan takes circumstances to be world-time pairs, CR expands the  
861 parameters to include all those parameters standardly recognized as elements  
862 of contexts of utterance (agents, times, places, etc.). Whereas circumstances  
863 of evaluation are usually appealed to in determining truth-value, CR allows  
864 them to determine content. Hence the same utterance can change its content  
865 (express a different proposition) if the context in which it is assessed changes.

866 Parsons (2011) considers the possibility of appealing to CR as a way of  
867 providing a semantic theory for answerphone cases. On the surface, the sug-  
868 gestion is promising: the utterance of “I am not here now” on the answerphone  
869 strikes us as intuitively true, despite the fact that it cannot be true if that sen-  
870 tence is contradictory. But, of course, CR will abandon the claim that it is  
871 contradictory, because that claim relies on the belief that the content of the  
872 sentence is tied to a proper context of utterance. CR can agree with Kaplan  
873 that all contexts of utterance are proper but take a more relaxed view on  
874 contexts of assessment, allowing these to impact the content of the sentence  
875 uttered in ways that break the tie with contexts of utterance. Hence we have a  
876 neat explanation of how an utterance of “I am not here now” can express a  
877 truth: although the sentence cannot be true when uttered, it can change its  
878 content depending on the context of its assessment so as to become true.

879 I think there is something right about CR, but we need to be careful about  
880 endorsing it as a semantic theory. The problem is that there is nothing sys-  
881 tematic about the behaviour of indexicals which tells us in advance whether  
882 they are assessment-sensitive or not. Answerphone messages are assessment-  
883 sensitive, ordinary utterances of indexical sentences tend not to be. Or, to be  
884 more precise, indexical sentences uttered in certain conventionally recognised  
885 scenarios are routinely interpreted in accordance with the predictions made  
886 by CR, while most utterances do not demand such elaborate mechanisms to  
887 interpret them. Of course, we might just maintain that CR applies uniformly  
888 to all utterances but that the default interpretation is one where the context  
889 of assessment coincides with the context of utterance. Only in certain cases  
890 does the context trigger a bifurcation of context of assessment from context  
891 of utterance. I see no problem with that view, but it clearly demonstrates  
892 that the semantic theory by itself does not do sufficient explanatory work. A  
893 pragmatic account of the way in which the interaction of context of utterance  
894 and context of assessment is triggered is essential to such a story, and this is  
895 what I have attempted to provide in this paper.



896 One thing that suggests that a CR-based semantics alone is not sufficient  
897 to explain displaced communication is that (as we have seen many times  
898 in the discussion in this paper) our intuitions are highly unpredictable and  
899 subject to the details of the contextual situation. Parsons takes this concern  
900 to show that CR cannot explain the answerphone problem. He imagines a  
901 case where a time delay on the phone line results in someone hearing the  
902 answerphone message after the speaker has in fact returned home. With some  
903 reservation, he endorses the view that the message is still true, and is (he  
904 claims) able to shift the context of assessment away from the time of decoding  
905 to the time of *intended* decoding. Parsons himself confesses to being unsure  
906 of his intuitions in regard to this example. It seems to me to be another case  
907 like those I considered previously which just show that we do not have firm  
908 intuitions about displaced communications. But without firm intuitions to  
909 make concrete predictions about what it meant and what is true or false, the  
910 task required of a semantic theory cannot be fulfilled. The situation can be  
911 seen quite clearly by reconsidering the holographic Ronnie scenario that I  
912 posed as an objection to Predelli's approach above. It is clear that using a  
913 holographic image of Ronnie requires some rich stage-setting to work. It is  
914 only because of this stage setting that the intended content (the pretence, as I  
915 have argued) is made available. We can recognize a semantic value which is  
916 interpreted relative to the context of assessment for the utterance. This will  
917 make sense of our intuition that some of Ronnie's apparent utterances at least  
918 sound like they are true ("it is raining in Manchester tonight," for example),  
919 while others don't sound true ("I am happy to be here tonight" doesn't sound  
920 true when we know that Real Ronnie is both dead and played no conscious  
921 role in this utterance). Consider Holographic Ronnie's production of "I am  
922 Ronnie James Dio"—is this true at the context of assessment? According to  
923 CR it ought to be possible that the agent really will be Real Ronnie. And Real  
924 Ronnie really is Ronnie James Dio. So the utterance should be true. But I  
925 don't have the intuition that this utterance is true—or, rather, I'm not sure  
926 that I have *any* intuition about this sort of case. Intuitions are just not stable  
927 in cases like these. And unstable intuitions are not suitable foundations for a  
928 semantic theory.

## 6 Objections and Replies

930 In this final section, I will consider some objections to the view that I have  
 931 presented above, and offer some replies, which will hopefully help to clarify  
 932 my position.

933 The first objection I want to consider concerns my definition of DCs. A  
 934 DC is a communication that occurs at a different context to that in which it  
 935 is encoded. It is tempting to assume (as seems to be the case for each of the  
 936 examples considered so far) that DCs are *always* evaluated with respect to  
 937 the context in which they are decoded (hence, on my view, the monstrosity  
 938 present in the pragmatic operation facilitating DCs) But what about cases  
 939 like we see in the following pair (both, imagine, recorded for a posthumously  
 940 broadcasted will):

941 (r) Today, I met with my lawyer before recording this will.

942 (s) Today, you all received a call from my lawyer informing you that you  
 943 have inherited a large sum.

944 It seems that both (r) and (s) are clear cases of DCs as commonly discussed in  
 945 the literature, yet only (s) seems to communicate information that is evaluated  
 946 with respect to the context in which it is decoded. Far from being monstrous,  
 947 (r) seems to communicate information about the context of encoding. But is  
 948 this not a DC?

949 I do not think that (r) is a DC. While (r) is being used to communicate  
 950 information at a context subsequent to that in which it is encoded, the infor-  
 951 mation is about the context of encoding. The indexicals “today,” “I,” “my,” and  
 952 “this” all contribute contents drawn from the context in which the message  
 953 is recorded. Furthermore, I am sceptical that a construction like (r) could be  
 954 developed in such a way as to be coherently understood as communicating  
 955 information about the context of decoding. For example, continuing (r) in the  
 956 following way, sounds infelicitous to my ear:

957 (r\*) Today, I met with my lawyer before recording this will that you are now  
 958 listening to.

959 If we understand this as an attempt to shift the temporal parameter of the  
 960 utterance from that indicated by “today” to that indicated by “now,” mid-  
 961 sentence, I think the sentence can only be made sense of if we read an implicit  
 962 open quotation as present on the “now.” Only in such cases, I suggest, do we

963 have a candidate for a DC. Simply presenting a recording of a message is not  
964 sufficient to produce a DC. Only when that message is naturally interpreted as  
965 communicating information *about* the context in which it is decoded, rather  
966 than encoded does it count as a DC. If I uncover a forgotten recording from  
967 my 10<sup>th</sup> birthday in which I say “I am 10 today,” I do not stumble on a DC. But  
968 if I uncover a recording of my 10-year-old self, saying “when you hear this,  
969 you will suddenly remember recording it when you were 10,” I do. It seems to  
970 me that (r) is akin to the former, not the latter.

971 The second objection arises when we consider a very large class of cases of  
972 potential DCs that I have said little about above, involving the production of  
973 *signs* containing indexicals. Consider a sign positioned in a hospital waiting  
974 room that consists of an inscription “please wait here.” This sign exhibits  
975 typical features of a DC as “here” will be naturally interpreted as referring to  
976 the location of installation, not of inscription. Tokens of the type of this sign  
977 are mass-produced in a factory. Some individual factory worker produced this  
978 particular token sign. But, surely, the producer of the sign in this case is not  
979 the agent of any instruction. The factory worker is simply a component part  
980 in the production of a communication that intuitively seems to occur at the  
981 time of decoding. This potentially casts doubt on my claim that DCs are the  
982 result of MOs operating on an utterance or inscription evaluated with respect  
983 to a proper Kaplanian context.<sup>19</sup>

984 I agree that it is implausible to construe the factory worker as the agent of  
985 the instruction inscribed on the sign. I am, however, unconvinced that *any*  
986 instruction as such is made in the factory. We should not be misled by the fact  
987 that human agents can be involved in the production of an artefact that carries  
988 information into inferring that they are the agents of whatever information  
989 is thereby transmitted. In this instance, the factory worker is no more the  
990 agent of a communicated content, than Stephen Hawkins’s voice-synthesising  
991 computer is the agent of his utterances when he relies on it to communicate  
992 his thoughts. The factory worker is producing another agent’s message in  
993 accordance with their instructions. Who, then, is the agent who desires to  
994 communicate the information? The agent here is the hospital (or relevant

---

19 Examples such as these motivate both O’Madagain (2014) and Briciu (2018) to distinguish between tokens and proper utterances. This distinction allows for the possibility of utterances at a distance by holding that genuine utterance requires the presence of illocutionary force, whereas the mere production of a token does not. I am inclined to agree that, in the example above, our factory worker is engaged in the production of a token, not an utterance but, as I now argue, I do not think that this means we must recognize the context of decoding as the context of utterance.

995 hospital authority). We can avoid complicated metaphysical questions about  
 996 how organisations might be agents by assuming an individual consultant,  
 997 Ms Smith, is the relevant authority. Ms Smith wants to ensure that patients  
 998 arriving in reception wait in an orderly fashion in the waiting room. One way  
 999 that she could do this would be to write a sign in her own hand saying “please  
 1000 wait here” ready to be displayed at the waiting room, or utter the sentence  
 1001 “please wait here” into a recording device to be on looped playback in the  
 1002 reception. But, due to the frequent reoccurrence of episodes when consultants  
 1003 need to instruct patients to wait in a specific location, it is of course more  
 1004 practical for signs to be mass-produced rather than produced by flimsy hand-  
 1005 written notes. Hence she orders a batch of ready-made signs designed to meet  
 1006 this common need among consultants. Nonetheless, Ms Smith remains the  
 1007 agent of the instruction. She has simply exploited a labour-saving device that  
 1008 ensures that one factory worker produces signs for the large number of agents  
 1009 who want to issue this instruction. Once she is in possession of the sign, she  
 1010 can exploit the convention that signs routinely signal information about their  
 1011 spatial location to engage in pretence of the sort her hand-written note would  
 1012 exploit.<sup>20</sup> The difference in the method of production of her message does not  
 1013 alter the fact that she is the agent of the utterance and its displacement is the  
 1014 result of a metalinguistic pretence, not a deferral of her utterance.

1015 Another objection responds directly to my analysis of open quotation as an  
 1016 MO. An obvious feature of open quotation is that, even if the quotation operation  
 1017 is not explicit, it should be easily recoverable. Consider this exchange  
 1018 from the movie *The Empire Strikes Back*. Lando Calrissian has double-crossed  
 1019 Han Solo and his friends, betraying them to Darth Vader and the evil Empire.  
 1020 However, he strikes a deal with Vader to preserve the freedom of Solo’s friends.  
 1021 Informing Solo of the deal, he says “I’ve done all I can. I’m sorry I can’t do  
 1022 more, but I’ve got my own problems.” Solo sarcastically replies: “Yeah. You’re  
 1023 a real hero.” It is obvious enough how we might appeal to an open quotation  
 1024 analysis of this ironic utterance. Solo is not expressing his own admiration for

---

20 Not all signs are obviously about the location in which they are placed, or object they are attached to, of course. An object may well be emblazoned with the sign “visit [such and such website] to see full product range,” or a pair of running shoes may come in a box marked “consult medical professional before beginning any new program of exercise.” Such signs, while clearly connected to some salient object are not about that object. But there are clearly a multitude of cases where the convention does hold: “twist clockwise” on a food jar lid, “made in England” on a guitar amplifier, “serve chilled” on a beer bottle, “4m high” on a road bridge, etc., all refer to the object they are attached to. “No smoking” in a public building, “Slow Down” on a road sign, “Wear a face covering” outside a shop, etc., all refer to the location in which they are placed.

1025 Calrissian (he has in fact just punched Calrissian in the face, unequivocally  
1026 expressing his real attitude). He does not mean that Calrissian is a hero, rather  
1027 he is mockingly echoing the use of this term of praise to display his own  
1028 distance from such a perspective. One thing that obviously stands in favour of  
1029 the open quotation analysis is that the recovery of the operation as an explicit  
1030 one is natural. One might very well report Solo's speech as "Yeah, You're a  
1031 real 'hero'." Many of the commonly cited cases of DCs in the literature on  
1032 indexicality, however, do not seem to be so neatly reconfigured with explicit  
1033 quotation marks. Consider Predelli's note from section 1:

1034 I am not at home now. If you hurry, you'll catch the evening flight  
1035 to Los Cabos. Meet me in six hours at the Hotel Cabo Real.

1036 It would not be natural to add quotation marks to the shifted indexicals in  
1037 the note (I add the "#" to indicate the marked quality of this):

1038 #I am not at home "now." If you hurry, you'll catch the evening  
1039 flight to Los Cabos. Meet me "in six hours" at the Hotel Cabo Real.

1040 The note, if anything, becomes quite confusing once the quotation marks  
1041 are made explicit. Why is this, and how can it be the case if a DC is really  
1042 generated by open quotation, in the same way as Han's response to Lando?

1043 A key difference is apparent in these two contrasting cases that explains  
1044 why quotation is not recoverable in the second case. The first case involves  
1045 the shifting of a sub-sentential element within a context that remains non-  
1046 shifted. The second case involves the shifting of the entire sentence for its  
1047 interpretation. But open quotation is linguistically employed for the first kind  
1048 of case only. Recall that open quotation is appealed to on my account as a  
1049 way of making explicit a form of mimicry. In a case where a single expression,  
1050 or string of expressions, contained in a wider linguistic frame are employed  
1051 in this mimicking role while the wider frame is not, open quotation serves  
1052 to explicitly indicate this role. When it is an entire sentence or other self-  
1053 contained linguistic item, this device serves no purpose. Mimicry shifts the  
1054 context to create a DC. Only when mimicry is embedded within a non-shifted  
1055 context is explicit quotation required to indicate this. This is why, for example,  
1056 it makes no sense to add quotation marks to this message:

1057 (t) "I'll be back"

1058 But they are clearly useful in:

1059 (t\*) As Arnold Schwarzenegger would say, “I’ll be back.”

1060 In Predelli’s note, we cannot capture the mimicry that I take it to employ by  
 1061 only quoting “now” and “in six hours,” because they have not been shifted in  
 1062 relation to the remainder of the note. It is the whole note that mimics what  
 1063 the note writer intends to communicate by their pretence. My claim is that  
 1064 the same sort of pretence lies behind both cases, but only when the mimicry  
 1065 is embedded within a wider frame is the quotation device demanded to make  
 1066 this explicit.

1067 The final objection I want to consider is a methodological one that, I think,  
 1068 goes to the heart of the different approaches to this problem taken by Predelli  
 1069 and those who, like myself, have urged a pragmatic explanation of DCs. The  
 1070 intuition that supports the pragmatic explanation of DCs is that distinctive  
 1071 features of the concrete episodes of language use that generate DCs are re-  
 1072 sponsible for them. These features point to ways in which DCs are heavily  
 1073 reliant on a rich background of human behaviour that goes beyond the remit  
 1074 of semantic theory to explain. Like other aspects of communication that are  
 1075 accepted as requiring an explanation from pragmatics such as conversational  
 1076 implicatures, we need to look beyond the literal meanings of the expressions  
 1077 employed to understand what is happening in these cases. But, insists Pre-  
 1078 delli (2005, 2011), this approach both misunderstands and underestimates  
 1079 the place of semantics. It fails, in his view, to respect Kaplan’s advice that we  
 1080 base our semantic theorizing on “the verities of meanings,” not “the vagaries  
 1081 of actions” (1989b, 585). Indeed, Kaplan’s own insistence on admitting only  
 1082 proper contexts is a failure to follow his own advice, according to Predelli  
 1083 (2005, 60–62). The view that Kaplan’s restriction on contexts strays from is  
 1084 summarised elsewhere by Predelli (2011, 301) like this: “[S]emantics is con-  
 1085 cerned with the evaluation of sentences with respect to contexts, not with  
 1086 concrete episodes of language use—it is concerned with ‘utterances’ only in  
 1087 Kaplan’s technical sense of the term as sentence-context pairs”. To illustrate  
 1088 the significance of drawing this methodological line, Predelli gives the ex-  
 1089 ample of a tautology that is sufficiently long that no mortal human will ever  
 1090 utter it. As a concrete episode of language use, we do not have an utterance.  
 1091 But clearly it is unproblematic to evaluate the sentence as true with respect to  
 1092 any context of “utterance,” in the more careful Kaplanian sense of a sentence-  
 1093 context pair (see 1989a, 522–523). Failing to respect this distinction, and being

1094 misled by the peculiarities of how we use language to perform speech acts  
1095 in particular situations, Predelli maintains, inevitably leads us to the wrong  
1096 conclusions.

1097 While I agree with Predelli that we ought to recognise the distinction he  
1098 makes, I do not agree that there is a methodological decision to be taken here  
1099 that will remain neutral with regard to our intuitions concerning concrete  
1100 episodes of language use. Keeping the discussion focused purely on the issue  
1101 at hand, one can of course construct a formal semantic theory that is more  
1102 generous in the contexts it evaluates sentences with respect to than the proper  
1103 contexts endorsed by Kaplan. One is limited only by mathematical constraints  
1104 in this regard. But eventually one has to make a decision about which of those  
1105 mathematical possibilities correspond to our *actual* use of natural language  
1106 expressions, if the mathematical structure in our formal semantics is going  
1107 to be empirically adequate as a model of the semantic profile of an actual  
1108 expression or set of expressions in a natural language. Kaplan's decision to  
1109 restrict the range of contexts we should be interested in to proper ones is, I  
1110 take it, based on this desideratum. After noting that an unconstrained range of  
1111 contexts will provide contexts with respect to which "I am here now" is false,  
1112 he insists that only the proper ones should be admitted if we are to arrive at  
1113 an empirically adequate analysis of the indexical expressions contained in  
1114 this sentence. To repeat the quotation I gave at the beginning of this paper:  
1115 "[I]mproper indices are like impossible worlds; no such contexts could exist  
1116 and thus there is no *interest* in evaluating the extensions of expressions with  
1117 respect to them" (Kaplan 1989a, 509, emphasis added). As the quotation shows,  
1118 while Predelli is quite correct to point out that a formal semantic theory should  
1119 be founded on an abstract pairing of expressions with mathematical objects  
1120 within a formal structure, we will have to make choices about which pairings  
1121 are of interest to our concerns as natural language semanticists, and these  
1122 choices will surely be based on our intuitions about the way the expressions  
1123 behave in the mouths, pens, and thoughts of ordinary speakers. My own  
1124 intuition, following Kaplan's, is that the restriction of contexts to the set of  
1125 proper contexts best captures the semantic behaviour of indexicals in English,  
1126 once we recognise the input of pragmatic processes on apparent deviations  
1127 from this restriction. Predelli's intuition is to take the deviations to illustrate  
1128 that Kaplan's restriction is empirically inadequate. I do not think that either  
1129 of us is basing our choice about which contexts our semantic theories should  
1130 recognise on issues independent of intuitions about concrete episodes of  
1131 language use, nor do I think that we should.

## 1132 7 Conclusion

1133 My argument in this paper has been the following. Firstly, I have argued that  
 1134 DCs are best understood as being generated by MOs. I have then argued that,  
 1135 understood as MOs, they are in turn best understood as the result of pragmat-  
 1136 ically triggered metalinguistic context-shifting operations. I have then given  
 1137 a detailed explanation of this proposed mechanism. If this is correct, then  
 1138 DCs are MOs but are not Monsters, for, while DCs are certainly monstrous,  
 1139 their monstrosity is not generated by any lexicalised semantic operator of En-  
 1140 glish. Furthermore, the argument presented here is also intended to vindicate  
 1141 Kaplan's insistence that the only proper contexts relevant to the semantic  
 1142 evaluation of English indexicals are those which situate the agent at the time  
 1143 and place of her utterance. DCs are not "utterances at a distance" which result  
 1144 from making utterances in improper contexts; they are ordinary utterances  
 1145 made in the course of a deliberate pretence that they are something more.\*

1146 Graham Stevens

1147  0000-0003-3832-6391

1148 University of Manchester

1149 Graham.P.Stevens@manchester.ac.uk

## 1150 References

- 1151 ÅKERMAN, Jonas. 2017. "Indexicals and Reference-Shifting: Towards a Pragmatic  
 1152 Approach." *Philosophy and Phenomenological Research* 95(1): 117–152, doi:10.111  
 1153 1/phpr.12216.
- 1154 BIANCHI, Claudia. 2014. "Slurs and Appropriation: An Echoic Account." *The Journal  
 1155 of Pragmatics* 66: 35–44, doi:10.1016/j.pragma.2014.02.009.
- 1156 BRICIU, Adrian. 2018. "Indexicals in Remote Utterances." *Philosophia* 46(1): 39–55,  
 1157 doi:10.1007/s11406-017-9909-x.
- 1158 CAPPELEN, Herman and LEPORE, Ernest. 1997. "The Varieties of Quotation." *Mind*  
 1159 106(424): 429–450, doi:10.1093/mind/106.423.429.
- 1160 CHUCK D and JAH, Yusuf. 1998. *Fight the Power: Rap, Race, and Reality*. New York:  
 1161 Dell Publishing.
- 1162 COHEN, Jonathan. 2013. "Indexicality and the Puzzle of the Answering Machine." *The  
 1163 Journal of Philosophy* 110(1): 5–32, doi:10.5840/jphil2013110143.

---

\* I would like to thank the audience at a seminar in the Linguistics department of the University of Manchester for their many helpful comments on a presentation of this paper. I am especially grateful to Justina Berškýtė for her many helpful comments on an earlier draft.



- 1164 CONNOLLY, Niall. 2017. "I Am Here Now, But I Won't Be Here When You Get This  
1165 Message." *Dialectica* 71(4): 603–622, doi:10.1111/1746-8361.12208.
- 1166 CORAZZA, Eros, FISH, William and GORVETT, Jonathan. 2002. "Who is 'I'?" *Philosophical  
1167 Studies* 107(1): 1–21, doi:10.1023/A:1013111419036.
- 1168 DAVIDSON, Donald. 1979. "Quotation." *Theory and Decision* 11(1): 27–40. Reprinted in  
1169 Davidson (1984, 479–488), doi:10.1007/BF00126690.
- 1170 —. 1984. *Inquiries into Truth and Interpretation*. Oxford: Oxford University Press,  
1171 doi:10.1093/0199246297.001.0001.
- 1172 EGAN, Andy. 2009. "Billboards, Bombs and Shotgun Weddings." *Synthese* 166(2):  
1173 251–279, doi:10.1007/s11229-007-9284-4.
- 1174 KAPLAN, David. 1989a. "Demonstratives." in *Themes from Kaplan*, edited by Joseph  
1175 ALMOG, John R. PERRY, and Howard K. WETTSTEIN, pp. 481–563. Oxford: Oxford  
1176 University Press. Widely circulated from 1977 on.
- 1177 —. 1989b. "Afterthoughts." in *Themes from Kaplan*, edited by Joseph ALMOG, John R.  
1178 PERRY, and Howard K. WETTSTEIN, pp. 565–614. Oxford: Oxford University Press.
- 1179 KING, Jeffrey C. 2001. *Complex Demonstratives: A Quantificational Account*. Cambridge,  
1180 Massachusetts: The MIT Press, doi:10.7551/mitpress/1990.001.0001.
- 1181 LASERSOHN, Peter. 2017. *Subjectivity and Perspective in Truth-Theoretic Semantics*.  
1182 Oxford Studies in Semantics and Pragmatics n. 8. Oxford: Oxford University Press,  
1183 doi:10.1093/acprof:oso/9780199573677.001.0001.
- 1184 MACFARLANE, John. 2014. *Assessment Sensitivity: Relative Truth and Its Applications*.  
1185 Oxford: Oxford University Press, doi:10.1093/acprof:oso/9780199682751.001.0001.
- 1186 O'MADAGAIN, Cathal. 2014. "Indexicals and the Metaphysics of Semantic Tokens:  
1187 When Shapes and Sounds become Utterances." *Thought* 3(1): 71–79, doi:10.1002/  
1188 tht3.114.
- 1189 PARSONS, Josh. 2011. "Assessment-Contextual Indexicals." *Australasian Journal of  
1190 Philosophy* 89(1): 1–17, doi:10.1080/00048400903493530.
- 1191 PREDELLI, Stefano. 1996. "Never Put Off Until Tomorrow What You Can Do Today."  
1192 *Analysis* 56(2): 85–91, doi:10.1093/analys/56.2.85.
- 1193 —. 2005. *Contexts: Meaning, Truth, and the Use of Language*. Oxford: Oxford University  
1194 Press, doi:10.1093/0199281734.001.0001.
- 1195 —. 2011. "I am Still Not Here Now." *Erkenntnis* 74(3): 289–303, doi:10.1007/s10670-  
1196 010-9224-4.
- 1197 —. 2014. "Kaplan's Three Monsters." *Analysis* 74(3): 389–393, doi:10.1093/analys/anu  
1198 059.
- 1199 RÉCANATI, François. 2010. *Truth-Conditional Pragmatics*. Oxford: Oxford University  
1200 Press, doi:10.1093/acprof:oso/9780199226993.001.0001.
- 1201 SCHLENKER, Philippe. 2003. "A Plea for Monsters." *Linguistics and Philosophy* 26(1):  
1202 29–120, doi:10.1023/A:1022225203544.
- 1203 SCOTT, Michael and STEVENS, Graham. 2019. "An Indexical Theory of Racial Pejora-  
1204 tives." *Analytic Philosophy* 60(4): 385–404, doi:10.1111/phib.12156.

- 1205 SHERMAN, Brett. 2015. "Constructing Contexts." *Ergo* 2(23): 581–605, doi:10.3998/er  
1206 go.12405314.0002.023.
- 1207 SIDELLE, Alan. 1991. "The Answering Machine Paradox." *Canadian Journal of Philos-*  
1208 *ophy* 21(4): 525–539, doi:10.1080/00455091.1991.10717260.
- 1209 STEVENS, Graham. 2009. "Utterance at a Distance." *Philosophical Studies* 143(2):  
1210 213–221, doi:10.1007/s11098-007-9199-4.
- 1211 STEVENS, Graham and DUCKETT, Nathan. 2019. "Expressive Content and Speaker-  
1212 Dependence." *Linguistic and Philosophical Investigations* 18: 97–112, doi:10.22381  
1213 /LPI1820195.
- 1214 VOLTOLINI, Alberto. 2006. "Fiction as a Base of Interpretation Contexts." *Synthese*  
1215 153(1): 23–47, doi:10.1007/s11229-006-0001-5.
- 1216 WEATHERSON, Brian. 2009. "Conditionals and Indexical Relativism." *Synthese* 166(2):  
1217 333–357, doi:10.1007/s11229-007-9283-5.
- 1218 WILSON, Deirdre. 2006. "The Pragmatics of Verbal Irony: Echo or Pretence?" *Lingua*  
1219 116(10): 1722–1743, doi:10.1016/j.lingua.2006.05.001.

# Constitutivism about Instrumental Desire and Introspective Belief

RYAN COX

This essay is about two familiar theses in the philosophy of mind: constitutivism about instrumental desires, and constitutivism about introspective beliefs, and the arguments for and against them. Constitutivism about instrumental desire is the thesis that instrumental desires are at least partly constituted by the desires and means-end beliefs which explain them, and is a thesis which has been championed most prominently by Michael Smith. Constitutivism about introspective belief is the thesis that introspective beliefs are at least partly constituted by the mental states they are about, and is a thesis which has been championed most prominently by Sydney Shoemaker. Despite their similarities, the fortunes of these two theses could not be more opposed: constitutivism about instrumental desire is widely accepted, and constitutivism about introspective belief is widely rejected. Yet, the arguments for both theses are roughly analogous. This essay explores these arguments. I argue that the argument which is widely taken to be the best argument for constitutivism about instrumental desires—what I call the argument from necessitation—does not provide the support for the thesis it is widely taken to provide, and that it fails for much the same reasons that it fails to provide support for constitutivism about introspective belief. Furthermore, I argue that the best argument for constitutivism about instrumental desires—what I will call the argument from cognitive dynamics—is also a good argument, if not equally good, for constitutivism about introspective belief (at least when the thesis is suitably qualified).

This essay is about two familiar theses in the philosophy of mind: *constitutivism about instrumental desire*, and *constitutivism about introspective belief*, and the arguments for and against them. Constitutivism about instrumental desire is the thesis that instrumental desires are at least partly constituted by the desires and means-end beliefs which explain them and is a thesis which has been championed most prominently by Michael Smith (2004). Consti-

1251 tutivism about introspective belief is the thesis that introspective beliefs are  
1252 at least partly constituted by the mental states they are about and is a the-  
1253 sis which has been championed most prominently by Sydney Shoemaker  
1254 (1996, 2012). Despite their similarities, the fortunes of these two theses could  
1255 not be more opposed: constitutivism about instrumental desire is widely ac-  
1256 cepted, and constitutivism about introspective belief is widely rejected. Yet,  
1257 the arguments for both theses are roughly analogous. So if one thesis is to be  
1258 accepted while the other is rejected there must be good reasons for rejecting  
1259 the arguments for one thesis but not the other.

1260 In this essay, I argue that the *best* argument for constitutivism about in-  
1261 strumental desire—what I will call the *argument from cognitive dynamics*—is  
1262 also a *good* argument, if not an *equally good* argument, for constitutivism  
1263 about introspective belief (at least when the thesis is suitably qualified). So,  
1264 at least with respect to this argument, there are no good reasons for accept-  
1265 ing one thesis while rejecting the other. At the same time, however, I argue  
1266 that the argument which is widely taken to be the best argument for con-  
1267 stitutivism about instrumental desire—what I will call the *argument from*  
1268 *necessitation*—does not provide the support for the thesis it is widely taken  
1269 to provide, and that it fails for much the same reasons that it fails to provide  
1270 support for constitutivism about introspective belief. So, with respect to this  
1271 argument, there are no good reasons for accepting one thesis while rejecting  
1272 the other, because this argument does not give us good reasons for accepting  
1273 either thesis.

1274 These conclusions suggest that the fortunes of constitutivism about instru-  
1275 mental desire and constitutivism about introspective belief are more closely  
1276 tied together than is often appreciated. Indeed, I hope to bring to bear on the  
1277 topic of instrumental desire an important lesson which has been learnt in the  
1278 philosophy of introspection. For philosophers of introspection have shown  
1279 that the argument from necessitation for introspective belief is unsuccess-  
1280 ful. This partly explains why constitutivism about introspective belief is not  
1281 widely accepted in the way constitutivism about instrumental desire is. Yet an  
1282 analogous lesson for the argument from necessitation for constitutivism about  
1283 instrumental desire has not yet been absorbed by those working on the philo-  
1284 sophy of instrumental desire. There is a certain irony here, since Shoemaker  
1285 explicitly draws analogies with the case of instrumental desire—particularly  
1286 the role played by means-end beliefs and non-instrumental desires in ratio-  
1287 nalisating and explaining further desires—in developing his arguments for  
1288 constitutivism about introspective belief (1996). Rejecting this argument for

1289 constitutivism about instrumental desire puts this thesis on a less firm footing,  
1290 and brings its fortunes more closely into line with those of constitutivism  
1291 about introspective belief. As I will argue, the fortunes of both theses rest  
1292 on the prospects of the argument from cognitive dynamics. I hope to show  
1293 that, even if the argument from cognitive dynamics is taken to offer adequate  
1294 support for constitutivism about instrumental desire while ultimately not  
1295 offering adequate support for constitutivism about introspective belief, the  
1296 fortunes of these two theses are tied closer together than is often appreciated.

1297 The essay is structured as follows. Section 1 discusses some preliminary  
1298 issues concerning how constitutivism about instrumental desire and consti-  
1299 tutivism about introspective belief are to be understood. Section 2 considers  
1300 the argument from necessitation for constitutivism about instrumental desire  
1301 and the argument from necessitation for constitutivism about introspective  
1302 belief respectively. I argue that both arguments fail. Section 3 considers the  
1303 argument from cognitive dynamics for constitutivism about instrumental  
1304 desire and the argument from cognitive dynamics for constitutivism about  
1305 introspective belief respectively. I argue that the arguments provide equally  
1306 good reasons for accepting both theses.

## 1307 **1 Locating the Topic**

1308 At the outset, we are going to need a way of understanding instrumental  
1309 desires and introspective beliefs which does not prejudice the case either for  
1310 or against the respective constitutivist theses. It might be thought, after all,  
1311 that there isn't much to say about constitutivism about instrumental desire,  
1312 since instrumental desires just are, by definition, according to some, those  
1313 desires which are at least partly constituted by the desires and means-end  
1314 beliefs which explain them. While the term "instrumental desire" is sometimes  
1315 used this way, there is a more neutral way of understanding instrumental  
1316 desires which does not prejudice the issue in this way. We can simply say that  
1317 instrumental desires just are, by definition, those desires which are *rationaly*  
1318 *explained* by other desires and means-end beliefs, where we leave open the  
1319 question of whether the former is partly constituted by the latter. This neutral  
1320 understanding of instrumental desire is widespread in the literature (Marks  
1321 1986, 9; Davis 1986, 69; Schroeder 2004, 5; McDaniel and Bradley 2008, 286;  
1322 Arpaly and Schroeder 2014, 6). It is arguably this understanding which Hume  
1323 has in mind when he writes:

1324 Ask a man *why he uses exercise*; he will answer *because he desires*  
 1325 *to keep his health*. If you then enquire, *why he desires to keep his*  
 1326 *health*, he will readily reply, *because his sickness is painful*. If you  
 1327 push your enquiries farther, and desire a reason *why he hates pain*,  
 1328 it is impossible he can ever give any. This is an ultimate end, and  
 1329 is never referred to any other object. (Hume EPM, Appendix I)<sup>1</sup>

1330 In this way, Hume distinguishes between what we would call the instrumen-  
 1331 tal desire to exercise, rationally explained by a desire to keep one's health,  
 1332 and the non-instrumental desire to avoid pain, not explained by any further  
 1333 desire. When Smith writes "I will call the desires which are explained by [non-  
 1334 instrumental] desires and means-end beliefs 'instrumental' desires" (2004,  
 1335 95) he is clearly stating that he means to use "instrumental desires" in this  
 1336 neutral way.<sup>2</sup>

1337 On this way of understanding instrumental desires, it is clearly a substantive  
 1338 question whether instrumental desires are partly constituted by means-end  
 1339 beliefs and other desires. Of course, many theorists accept this substantive  
 1340 thesis about instrumental desires. Here is Smith's assertion of his commitment  
 1341 to constitutivism about instrumental desire:

1342 Instrumental desires are not distinct from the non-instrumental  
 1343 desires and means-end beliefs that explain them, but are rather  
 1344 just the complex state of having such non-instrumental desires  
 1345 and means-end beliefs standing in a suitable relation. (2004, 96)

1346 According to Smith, then, instrumental desires are (at least partly, if not  
 1347 wholly) constituted by non-instrumental desires and means-end beliefs stand-  
 1348 ing in suitable relations. As he puts it: "Instrumental desires are thus better  
 1349 thought of as being nothing over and above the non-instrumental desires  
 1350 and means-end beliefs that explain them" (2004, 96).<sup>3</sup> Here Smith appears to  
 1351 commit himself to a strong form of constitutivism about instrumental desire,  
 1352 one which holds that instrumental desires are "nothing over and above" or are  
 1353 *wholly constituted* by the non-instrumental desires and means-end beliefs that

---

1 Quoted in Smith (2004, 94).

2 See also "[I]nstrumental desires are those that can be explained by non-instrumental desires and means-end beliefs" (Smith 2004, 96–97).

3 The qualification that the non-instrumental desire and the means-end belief must stand in a suitable relation is missing here but we can assume that it is intended.

1354 explain them standing in a suitable relation. Constitutivism about instrumen-  
1355 tal desire, so understood, is a widely accepted thesis in the philosophy of mind.  
1356 However, as we have just seen, it is a substantive thesis about instrumental  
1357 desires, and we will be concerned with the arguments for it in the next two  
1358 sections.

1359 Introspective beliefs raise their own unique problems for understanding.  
1360 While there is no temptation in the case of introspective beliefs to hold that  
1361 they are, by definition, beliefs which are partly constituted by the mental  
1362 states they are about, it is nonetheless difficult to say what an introspective  
1363 belief is. It is not enough to say that an introspective belief is just a belief  
1364 about one's own mental states. For, it is relatively uncontroversial that at  
1365 least some of our beliefs about our own mental states are arrived at on the  
1366 basis of inference. And these beliefs are very plausibly thought of as being  
1367 constitutively distinct from the mental states they are about. They are no more  
1368 constituted by the mental states they are about than our beliefs about the  
1369 mental states of others are constituted by the mental states they are about.  
1370 Rather, introspective beliefs must be understood as beliefs about our own  
1371 mental states which are arrived at by some special means and are not based  
1372 on evidence or observation in the way that our beliefs about the mental states  
1373 of others are.<sup>4</sup> While it is a matter of controversy how introspective beliefs  
1374 are to be understood, I will assume that some such distinction can be drawn  
1375 among beliefs about our own mental states and that only some of these beliefs  
1376 will count as introspective beliefs. Constitutivism about introspective belief  
1377 holds that *these* beliefs are at least partly constituted by the mental states they  
1378 are about. Here is Shoemaker's assertion of his commitment to constitutivism  
1379 about introspective belief:

1380       What I am inclined to say is that second-order belief, and the  
1381       knowledge it typically embodies, is supervenient on first order  
1382       beliefs and desires—or rather, it is supervenient on these plus a  
1383       certain degree of rationality, intelligence, and conceptual capacity.  
1384       By this I mean that one has the former *in* having the latter—that  
1385       having the former is nothing over and above having the latter.  
1386       (1996, 34)

---

4 For an influential discussion of these features of introspective beliefs see Moran (2001). For scepticism about the existence of introspective beliefs so understood see Cassam (2014).

1387 According to Shoemaker, then, introspective beliefs are at least partly (if not  
 1388 wholly) constituted by the mental states they are about together with “a cer-  
 1389 tain degree of rationality, intelligence, and conceptual capacity”—having the  
 1390 former is “nothing over and above” having the latter. This passage anticipates  
 1391 the argument to be examined in the next section, as Shoemaker moves here  
 1392 from a claim about supervenience or necessitation to a claim about consti-  
 1393 tution. So we have seen that constitutivism about introspective belief is a  
 1394 substantive thesis about introspective belief. It will be the aim of the next  
 1395 sections of the essay to evaluate arguments for both forms of constitutivism.

1396 Before turning to those arguments a final qualification is in order. In the  
 1397 passage from Shoemaker just quoted, Shoemaker is concerned with a kind of  
 1398 constitutivism about introspective belief which concerns introspective beliefs  
 1399 about *attitudes* in particular. For the most part, in what follows, I will be  
 1400 concerned with forms of constitutivism about introspective belief which are  
 1401 restricted in this way, holding that introspective beliefs about our *attitudes*  
 1402 are partly constituted by the attitudes they are about. Introspective beliefs  
 1403 about phenomenally conscious states raise further issues that I will not be  
 1404 able to address here, and while constitutivism about these introspective beliefs  
 1405 may be defensible, different arguments may be required.<sup>5</sup> There is also an  
 1406 interesting question about how constitutivism might be extended to states  
 1407 which are plausibly thought of as having both cognitive and non-cognitive  
 1408 components, like emotions. While I think that constitutivism can be defended  
 1409 for a wide range of mental states, I will largely set aside such an exploration  
 1410 here, and will focus on the particular case of introspective beliefs about our  
 1411 own attitudes. I will also set aside the difficult question of whether we can or  
 1412 should expect a uniform account of introspection and introspective beliefs  
 1413 which applies to all mental states.<sup>6</sup> Finally, related to these questions is the  
 1414 question of which theories of introspection—theories of the means by which  
 1415 we arrive at introspective beliefs—are compatible with constitutivism about  
 1416 introspective belief and those which are not. Some theories of introspection  
 1417 have implications for the relation between introspective beliefs and the mental  
 1418 states they are about. The self-scanning theory of David Armstrong (1968)  
 1419 has the implication that introspective beliefs are constitutively distinct from  
 1420 the states they are about since the former are caused by the latter. Other  
 1421 theories are neutral about the relation and while they may provide causal

---

5 See Chalmers (2010) for a discussion of phenomenal beliefs.

6 See Boyle (2009) and Byrne (2011) for discussion.



1422 explanations of the means by which we arrive at introspective beliefs—by  
1423 means of answering deliberative questions (Moran 2001), or by means of an  
1424 ascent-routine (Gordon 1986)—they remain compatible with the possibility  
1425 that introspective beliefs are partly constituted by the attitudes they are about.  
1426 I will also set aside this question, as the arguments we will consider for  
1427 constitutivism about introspective belief do not presuppose any particular  
1428 theory of introspection.

## 1429 2 The Arguments from Necessitation

1430 In this section I will formulate and evaluate the argument from necessitation  
1431 for constitutivism about instrumental desire and the argument from neces-  
1432 sitation for constitutivism about introspective belief. The arguments belong  
1433 to a family of arguments, arguments from necessitation or supervenience,  
1434 which are familiar enough across many areas of philosophy. The crucial step  
1435 in such arguments is a move from a claim about metaphysical necessitation  
1436 or supervenience, to a claim about constitution. I will first formulate and  
1437 motivate each argument and then turn to evaluation.

### 1438 2.1 For Constitutivism About Instrumental Desire

1439 For simplicity, we can formulate and motivate the argument from necessita-  
1440 tion for constitutivism about instrumental desire by focussing on an arbitrary  
1441 example which we can take to reveal something general about instrumental  
1442 desires. Suppose, then, that Jane desires to exercise because she desires to keep  
1443 her health, and believes that exercising is a means to keeping her health. It  
1444 follows from our understanding of instrumental desires, and the assumption  
1445 that this is the “because” of rational explanation, that Jane instrumentally  
1446 desires to exercise. Why think that Jane’s desiring to exercise is at least partly  
1447 constituted by her desiring to keep her health and her believing that exercising  
1448 is a means to keeping her health? The argument from necessitation proceeds  
1449 in two steps.

1450 The first step establishes that the relation between Jane’s desiring to keep  
1451 her health, her believing that exercising is a means to keeping her health, *her*  
1452 *being fully rational*, and her desiring to exercise is not a merely contingent  
1453 relation, but is, in some sense, necessary. We can bring this out by reflecting  
1454 on a claim about necessity like the following:

1455 N<sub>1</sub>. Necessarily, if Jane desires to keep her health, believes that  
 1456 exercising is a means to keeping her health, and is fully rational,  
 1457 then she desires to exercise.

1458 I will have more to say about the relevant understanding of “is fully rational”  
 1459 in this claim below, and will connect this to Smith’s claims about means-  
 1460 end beliefs and non-instrumental desires standing in suitable relations. But  
 1461 for now, we can simply observe that on a natural understanding of “is fully  
 1462 rational” such a claim is intuitively plausible. To bring this out we might  
 1463 notice that while it is certainly possible for Jane to desire to keep her health,  
 1464 to believe that exercising is a means to keeping her health, while not desiring  
 1465 to exercise to any degree—after all, she might be less than fully rational—it is  
 1466 not possible for her to desire to keep her health, to believe that exercising is a  
 1467 means to keeping her health, and not desire to exercise to any degree *if she is*  
 1468 *fully rational*.

1469 The second step establishes the best explanation of the necessary connec-  
 1470 tion is that Jane’s desiring to exercise is at least partly, if not wholly, constituted  
 1471 by her desiring to keep her health, her believing that exercising is a means to  
 1472 keeping her health, and her being fully rational. While it might be tempting  
 1473 to move directly from the claim about necessitation to this conclusion, it  
 1474 is important to see that there are alternative explanations of the necessary  
 1475 connection which need to be considered. The explanations of the necessary  
 1476 connection we need to consider in this case are these:

1477 H<sub>1</sub>. Jane’s desiring to exercise is wholly constituted by her desiring to  
 1478 keep her health, her believing that exercising is a means to keeping  
 1479 her health, and her being fully rational.

1480 H<sub>2</sub>. Jane’s desiring to keep her health, her believing that exercising  
 1481 is a means to keeping her health, and her being fully rational is  
 1482 partly constituted by her desiring to exercise.

1483 H<sub>3</sub>. There is something that Jane’s desiring to exercise is wholly  
 1484 constituted by and which Jane’s desiring to keep her health, her  
 1485 believing that exercising is a means to keeping her health and her  
 1486 being fully rational is partly constituted by.

1487 These explanations correspond to the familiar options for explaining nec-  
1488 cessary connections of this form. If there is a necessary connection between  
1489 something's being an *F* and its being both a *G* and an *H*, then we might explain  
1490 this necessary connection by holding that its being an *F* is wholly constituted  
1491 by its being a *G* and an *H*, or by holding that its being a *G* and an *H* is partly  
1492 constituted by its being an *F*, or by holding that there is something else, such  
1493 that its being a *G* and an *H* is partly constituted by and its being an *F* is wholly  
1494 constituted by.

1495 Of these explanations,  $H_1$  certainly looks to be the best. There is no candi-  
1496 date for the kind of third-factor required by  $H_3$ , and, at least initially, it is hard  
1497 to see why Jane's desiring to keep her health, her believing that exercising is  
1498 a means to keeping her health, and her being fully rational, would be partly  
1499 constituted by her desiring to exercise: none of these conditions seem to be  
1500 independently partly constituted by her desiring to exercise, and it is hard to  
1501 see how, jointly, they could be partly constituted by her desiring to exercise.  
1502 So we may tentatively conclude that  $H_1$  provides the best explanation of the  
1503 necessary connection.

1504 These two steps, then, provide motivation for the two premises of the  
1505 argument from necessitation for constitutivism about instrumental desire.  
1506 We can think of the argument as proceeding as follows:

- 1507  $P_1$  There is a necessary connection between (i) Jane's desiring to exercise  
1508 and (ii) her desiring to keep her health, her believing that exercising is  
1509 a means to keeping her health, and her being fully rational.
- 1510  $P_2$  The best explanation of this necessary connection is that Jane's desiring  
1511 to exercise is wholly constituted by her desiring to keep her health, her  
1512 believing that exercising is a means to keeping her health, and her being  
1513 fully rational.
- 1514  $C_1$  Jane's desiring to exercise is wholly constituted by her desiring to keep  
1515 her health, her believing that exercising is a means to keeping her health,  
1516 and her being fully rational.

1517 Since the example we have focused on here was entirely arbitrary, the same  
1518 reasoning can be followed in arguing for constitutivism about instrumental  
1519 desire as a general thesis. At least initially, then, the argument from necessita-  
1520 tion provides a good case for constitutivism about instrumental desire. I will  
1521 now explicate an analogous argument for constitutivism about introspective  
1522 belief before turning to objections.

## 2.2 For Constitutivism About Introspective Belief

1524 For simplicity, we can formulate and motivate the argument from necessita-  
 1525 tion for constitutivism about introspective belief by focusing on an arbitrary  
 1526 example which we can take to reveal something general about introspective  
 1527 beliefs. Suppose that Jane introspectively believes that she believes it is about  
 1528 to rain. Why think that Jane's believing that she believes it is about to rain  
 1529 is not constitutively distinct from her believing that it is about to rain? The  
 1530 argument proceeds in two steps.

1531 The first step establishes that the relation between Jane's believing that she  
 1532 believes that it is about to rain and her believing that it is about to rain, her  
 1533 having some interest in the question of whether she believes that it is about  
 1534 to rain, and her being fully rational, is not a merely contingent relation, but  
 1535 is, in some sense, necessary. We can bring this out by reflecting on a claim  
 1536 about necessity like the following:

1537  $N_2$ . Necessarily, if Jane believes that it is about to rain, understands  
 1538 and has some interest in the question of whether she believes that  
 1539 it is about to rain, and is fully rational, then she believes that she  
 1540 believes that it is about to rain.

1541 A few clarifications are in order here. Recall Shoemaker's claim that introspec-  
 1542 tive beliefs are supervenient on "a certain degree of rationality, intelligence,  
 1543 and conceptual capacity." This suggests a claim like the following: necessarily,  
 1544 if Jane believes that it is about to rain, has a certain degree of rationality,  
 1545 intelligence, and conceptual capacity, she believes that it is about to rain.  
 1546 Shoemaker adds these claims about intelligence and conceptual capacity here  
 1547 in order to avoid problems stemming from small children and animals who  
 1548 may have the relevant attitudes, have a certain degree of rationality, and yet  
 1549 not even be able to understand the question of whether they have the attitudes  
 1550 in question. I have captured this element of Shoemaker's view with the claim  
 1551 about understanding in  $N_2$ . However, I have added the further claim that Jane  
 1552 must *have some interest in the question of whether she believes that it is about to*  
 1553 *rain*. Arguably, Shoemaker's conditions are too weak. It seems to be possible  
 1554 for Jane to believe that it is about to rain, for her to understand the question  
 1555 of whether she believes that it is about to rain, for her to be fully rational,  
 1556 and yet for her not to believe that she believes that it is about to rain *if she*

1557 *has no interest in the question of whether she believes that it is about to rain.*<sup>7</sup>  
1558 Quite generally, it seems that we will not believe that we have some attitude  
1559 or another when we have no interest in the question of whether we have  
1560 that attitude. With these clarifications in order, we can see that this claim is  
1561 intuitively plausible on a natural understanding of “is fully rational.” To bring  
1562 this out we might notice that while it is certainly possible for Jane to believe  
1563 that it is about to rain, to have some interest in the question of whether she  
1564 believes that it is about to rain, while not believing that she believes that it is  
1565 about to rain—after all, she might be less than fully rational—it is not possible  
1566 for her to believe that it is about to rain, to have some interest in the question  
1567 of whether she believes that it is about to rain, and to not believe that she  
1568 believes that it is about to rain *if she is fully rational*.

1569 The second step in the argument establishes that the best explanation of  
1570 this necessary connection is that Jane’s believing that she believes that it  
1571 is about to rain is at least partly, if not wholly, constituted by her believing  
1572 that it is about to rain, her having some interest in the question of whether  
1573 she believes that it is about to rain, and her being fully rational. Again, we  
1574 must consider the alternative explanations. The explanations of the necessary  
1575 connections we need to consider in this case are these:

1576 H<sub>1</sub>. Jane’s believing that she believes that it is about to rain is wholly  
1577 constituted by her believing that it is about to rain, her understand-  
1578 ing and taking an interest in the question of whether she believes  
1579 that it is about to rain, and her being fully rational.

1580 H<sub>2</sub>. Jane’s believing that it is about to rain, her understanding and  
1581 taking an interest in the question of whether she believes that it is  
1582 about to rain, and her being fully rational is partly constituted by  
1583 her believing that she believes that it is about to rain.

1584 H<sub>3</sub>. There is something that Jane’s believing that she believes that it  
1585 is about to rain is wholly constituted by and which Jane’s believing  
1586 that it is about to rain, her understanding and having an interest in  
1587 the question of whether she believes that it is about to rain, and her  
1588 being fully rational is partly constituted by.

---

7 See Stoljar (2019) for a discussion of further ways of refining such claims. For the purposes of the argument, all that matters is that there is some non-trivial, finite, mental condition which necessitates the introspective belief.

1589 Of these explanations,  $H_1$  certainly looks to be the best. There is no candidate  
 1590 for the kind of third-factor required by  $H_3$ , and, at least initially, it is hard  
 1591 to see why Jane's believing that it is about to rain, her understanding and  
 1592 taking an interest in the question of whether she believes that it is about to  
 1593 rain, and her being fully rational, would be partly constituted by her believing  
 1594 that she believes that it is about to rain: none of these conditions seem to be  
 1595 independently partly constituted by her believing that she believes that it is  
 1596 about to rain, and it is hard to see how, jointly, they could be partly constituted  
 1597 by her believing that she believes that it is about to rain. So we may tentatively  
 1598 conclude that  $H_1$  provides the best explanation of the necessary connection.

1599 These two steps, then, provide motivation for the two premises of the  
 1600 argument from necessitation for constitutivism about introspective belief. We  
 1601 can think of the argument as proceeding as follows:

- 1602  $P_1$  There is a necessary connection between (i) Jane's believing that she  
 1603 believes that it is about to rain and (ii) her believing that it is about to  
 1604 rain, her understanding and having some interest in the question of  
 1605 whether she believes that it is about to rain, and her being fully rational.  
 1606  $P_2$  The best explanation of this necessary connection is that Jane's believing  
 1607 that she believes that it is about to rain is wholly constituted by her  
 1608 believing that it is about to rain, her understanding and having some  
 1609 interest in the question of whether she believes that it is about to rain,  
 1610 and her being fully rational.  
 1611  $C_1$  Jane's believing that she believes that it is about to rain is wholly consti-  
 1612 tuted by her believing that it is about to rain, her understanding and  
 1613 having some interest in the question of whether she believes that it is  
 1614 about to rain, and her being fully rational.

1615 Since the example we have focused on here was entirely arbitrary, the same  
 1616 reasoning can be followed in arguing for constitutivism about introspective be-  
 1617 lief as a general thesis. At least initially, then, the argument from necessitation  
 1618 provides a good case for constitutivism about introspective belief.

### 2.3 *Evaluating the Arguments*

1620 We can now turn to the evaluation of the arguments from necessitation.  
 1621 Whether the arguments are successful turns crucially on how the notion of  
 1622 rationality is understood. So far I have presented the arguments without com-

1623 ment on how rationality is to be understood. I will now argue that there are  
1624 two understandings of “being fully rational” which are relevant to the argu-  
1625 ments from necessitation: an evaluative sense and a dispositional sense. When  
1626 the arguments are understood in terms of the former, their first premises are  
1627 true, but their second premises are false. When the arguments are under-  
1628 stood in terms of the latter, their first premises are false. The arguments get  
1629 whatever force they have from equivocating on these two understandings of  
1630 “being fully rational.” This objection to the arguments from necessitation is  
1631 due, in its essentials, to Amy Kind, who makes the objection in connection  
1632 with the argument from necessitation for constitutivism about introspective  
1633 belief (2003).<sup>8</sup> While the objection has been generally appreciated in the philo-  
1634 sophical literature on introspection (Gertler 2010), it has not been generally  
1635 appreciated in the philosophical literature on instrumental desires.<sup>9</sup>

1636 It is natural to think about the arguments above on an evaluative under-  
1637 standing of “being fully rational.” On this understanding, someone is fully  
1638 rational if and only if they are not in violation of the principles of rationality,  
1639 that is, if and only if they fully conform to the principles of rationality. On this  
1640 understanding, someone is less than fully rational if they do not fully conform  
1641 to the principles of rationality. It is very plausible that there is a principle of  
1642 rationality which requires you to desire the means if you desire some end and  
1643 believe that the means are a means to that end. Similarly, it is very plausible  
1644 that there is a principle of rationality which requires you to believe that you  
1645 have some attitude if you have that attitude and you understand and have  
1646 some interest in the question of whether you have it. This can be brought out  
1647 by reflecting on the necessitation claims with this understanding of “being  
1648 fully rational” made fully explicit:

1649  $N_1'$ . Necessarily, if Jane desires to keep her health, believes that  
1650 exercising is a means to keeping her health, and fully conforms to  
1651 the principles of rationality, then she desires to exercise.

1652  $N_2'$ . Necessarily, if Jane believes that it is about to rain, understands  
1653 and has some interest in the question of whether she believes that

---

8 While the objection, in its essentials, is due to Kind, the specific development of the objection made here is original to this essay.

9 In both Kind's (2003) and Gertler's (2010) discussions, there is an appeal to causation, and a contrast between causation and constitution, which is not made in the presentation of the objection here.

1654 it is about to rain, and fully conforms to the principles of rationality,  
 1655 then she believes that she believes that it is about to rain.

1656 Since it is arguably this understanding of “being fully rational” which we  
 1657 evaluated the original necessitation claims with, it is not surprising that they  
 1658 both appear to be plausible when this understanding is made fully explicit.

1659 The problem arises for the arguments from necessitation when we turn  
 1660 to the evaluation of their second premises on this understanding. To see the  
 1661 problem in the case of instrumental desire, notice that if Jane desires to keep  
 1662 her health and believes that exercising is a means to keeping her health, then,  
 1663 in order to fully conform to the principles of rationality, she must desire to  
 1664 exercise. But then, given that she desires to keep her health and believes  
 1665 that exercising is a means to keeping her health, *if* she fully conforms to  
 1666 the principles of rationality, this must be at least partly because she desires  
 1667 to exercise. So, her fully conforming to the principles of rationality is partly  
 1668 constituted by her desiring to exercise. And if her fully conforming to the  
 1669 principles of rationality is partly constituted by her desiring to exercise, her  
 1670 desiring to exercise cannot be even partly constituted by her conforming to  
 1671 the principles of rationality. So, on this understanding of “being fully rational,”  
 1672  $H_1$  is not the best explanation of the necessary connection and the argument  
 1673 from necessitation fails.

1674 To see the problem in the case of introspective belief, notice that if Jane  
 1675 believes that it is about to rain, and understands and has some interest in the  
 1676 question of whether she believes that it is about to rain, then, in order to fully  
 1677 conform to the principles of rationality, she must believe that she believes that  
 1678 it is about to rain. But then, given that she believes that it is about to rain, and  
 1679 understands and has some interest in the question of whether she believes  
 1680 that it is about to rain, *if* she fully conforms to the principles of rationality,  
 1681 this must be at least partly because she believes that she believes that it is  
 1682 about to rain. So, her fully conforming to the principles of rationality is partly  
 1683 constituted by her believing that she believes that it is about to rain. And if  
 1684 her fully conforming to the principles of rationality is partly constituted by  
 1685 her believing that she believes that it is about to rain, her believing that she  
 1686 believes that it is about to rain cannot be even partly constituted by her fully  
 1687 conforming to the principles of rationality. So on this understanding  $H_1$  is  
 1688 not the best explanation of the necessary connection and the argument from  
 1689 necessitation fails.



1690 An analogy might help to drive home the crucial point here. Suppose that  
1691 you are legally required to pay your taxes before the end of the financial year  
1692 every financial year. Then, you fully conform to the law only if you pay your  
1693 taxes before the end of the financial year every year. Now, if it is the end of  
1694 the financial year, then *if* you fully conform to the law, this must be at least  
1695 partly because you have paid your taxes. So, your fully conforming to the law  
1696 is partly constituted by your paying your taxes. And if your fully conforming  
1697 to the law is partly constituted by your paying your taxes, then your paying  
1698 your taxes cannot be, even partly, constituted by your fully conforming to  
1699 the law. As this analogy demonstrates, while conforming to certain norms  
1700 may, along with other conditions, necessitate certain further conditions, it is  
1701 your conforming to the norms which is partly constituted by those further  
1702 conditions, and not vice versa. It is not surprising then that on the evaluative  
1703 understanding of “being fully rational” the second premises of the arguments  
1704 from necessitation are false.

1705 At this stage a proponent of the arguments from necessitation may argue  
1706 that this is not the understanding of “being fully rational” they had in mind.  
1707 They may instead appeal to a dispositional understanding of “being fully  
1708 rational.” On this understanding, being fully rational is being disposed to  
1709 conform to the principles of rationality (perhaps along with there being no  
1710 barrier to one’s manifesting this disposition). Someone who was fully disposed  
1711 to conform to the principles of rationality, where there is no barrier to their  
1712 manifesting this disposition, would very plausibly *come to desire* to exercise if  
1713 they desired to keep their health and believed that exercising was a means  
1714 to keeping their health. Someone who was fully disposed to conform to the  
1715 principles of rationality, where there was no barrier to their manifesting this  
1716 disposition, would very plausibly *come to believe* that they believed that it was  
1717 about to rain if they believed that it was about to rain and understood and  
1718 had some interest in the question of whether they believed that it was about  
1719 to rain. The trouble with this understanding is that the first premises of the  
1720 arguments from necessitation seem to be false if they are understood in terms  
1721 of it. To see the trouble, we can reflect on the necessitation claims with this  
1722 understanding of “being fully rational” made fully explicit:

1723  $N_1$  ". Necessarily, if Jane desires to keep her health, believes that exer-  
1724 cising is a means to keeping her health, is disposed to conform to the  
1725 principles of rationality, and there is no barrier to the manifestation  
1726 of her disposition, then she desires to exercise.

1727  $N_2''$ . Necessarily, if Jane believes that it is about to rain, understands  
 1728 and has some interest in the question of whether she believes that it  
 1729 is about to rain, is disposed to conform to the principles of rationality,  
 1730 and there is no barrier to the manifestation of her disposition, then  
 1731 she believes that she believes that it is about to rain.

1732 The problem here is that even when all barriers to the manifestation of a  
 1733 disposition are removed, there is no metaphysically necessary connection  
 1734 between having the disposition, the triggering conditions for the disposition  
 1735 obtaining, and the disposition manifesting. At best, the relation between  
 1736 having a disposition, the triggering conditions for that disposition obtaining,  
 1737 and the manifestation of the disposition, is one of *nomological necessity*. On  
 1738 this understanding of “being fully rational” the first premises of the arguments  
 1739 from necessitation are false, and there is no necessary connection between  
 1740 being in certain psychological conditions, being fully rational, and being in  
 1741 some further psychological condition.

1742 It might be thought at this point that we might add some further condition to  
 1743 the antecedents of the conditionals above to avoid this problem. My response  
 1744 to this suggestion, however, is that there is nothing that can be added that  
 1745 will not either (i) fall prey to considerations like those just given against the  
 1746 current proposal or (ii) fall prey to considerations like those given against  
 1747 the evaluative understanding of “being fully rational.” Suppose someone  
 1748 suggests that we need only add to the antecedents of the conditionals that the  
 1749 dispositions be manifested. Then we would have:

1750  $N_1'''$ . Necessarily, if Jane desires to keep her health, believes that  
 1751 exercising is a means to keeping her health, is disposed to conform  
 1752 to the principles of rationality, and manifests this disposition, then  
 1753 she desires to exercise.

1754  $N_2'''$ . Necessarily, if Jane believes that it is about to rain, and has  
 1755 some interest in the question of whether she believes that it is about  
 1756 to rain, is disposed to conform to the principles of rationality, and  
 1757 manifests this disposition, then she believes that she believes that it  
 1758 is about to rain.

1759 But this suggestion faces a problem analogous to that faced by the view which  
 1760 understands “being fully rational” in the evaluative sense. Given that Jane

1761 desires to keep her health, believes that exercising is a means to keeping  
1762 her health, and is disposed to conform to the principles of rationality, *if* she  
1763 manifests her disposition, this must be in part because she desires to exercise.  
1764 To manifest this disposition is, in part, to come to desire to exercise. So, her  
1765 manifesting this disposition is partly constituted by her desiring to exercise.  
1766 And if her manifesting this disposition is partly constituted by her desiring to  
1767 exercise, her desiring to exercise cannot be even partly constituted by her man-  
1768 ifesting this disposition (and so cannot be wholly constituted by her desiring  
1769 to keep her health, believing that exercising is a means to keeping her health,  
1770 being disposed to conform to the principles of rationality, and manifesting  
1771 this disposition). I suspect that Smith’s appeal to the idea that a means-end  
1772 belief and a non-instrumental desire must stand “in a suitable relation” in his  
1773 statement of constitutivism about instrumental desire is an attempt to straddle  
1774 the gap here between the triggering conditions of a disposition obtaining and  
1775 the disposition manifesting. As we have just seen, however, if we understand  
1776 the claim in terms of the triggering conditions being met—the means-end  
1777 belief and the non-instrumental desire are appropriately related just when  
1778 they trigger the relevant disposition—then Smith’s position falls prey to con-  
1779 siderations like those given against the dispositional understanding, and if  
1780 we understand it in terms of the disposition manifesting—the means-end  
1781 belief and the non-instrumental desire are appropriately related just when  
1782 the relevant disposition manifests—then it falls prey to the considerations  
1783 just given.

1784 As initially compelling as the arguments from necessitation may seem, they  
1785 are ultimately unsuccessful. As I said earlier, this is old news in the literature  
1786 on introspective belief. There, it is widely conceded that Sydney Shoemaker  
1787 may be right about the necessitation claim, at least if it is understood on the  
1788 evaluative sense of “rational,” but it is held that nothing follows from this  
1789 vis-à-vis the constitutive theory of introspective belief. The problem is that  
1790 Shoemaker moves too quickly from the necessitation claim to the constitution  
1791 claim. To see this, consider the passage quoted earlier from Shoemaker. There  
1792 Shoemaker moves from a supervenience claim in the first sentence to the  
1793 constitution claim in the second. Indeed, he says that he means the same  
1794 thing by both claims. But, at best, he has only argued for the metaphysical

1795 necessitation or supervenience claim, and the constitution claim does not  
1796 immediately follow, as we have seen.<sup>10</sup>

1797 It is doubtful that Shoemaker had the evaluative sense of “rational” in mind,  
1798 however. In one of the only places where Shoemaker gives us any clues about  
1799 the sense of rationality he has in mind, he writes: “The fact that the person is  
1800 rational might be compared to the fact that the powder in the bomb was dry”  
1801 (1996, 32). This certainly suggests that Shoemaker had the dispositional sense  
1802 in mind. And, as we have seen, it is plausible that *if* the necessitation claim  
1803 were true on the dispositional sense of “rational,” then the constitution claim  
1804 would follow. It is widely agreed, however, that Shoemaker’s arguments for  
1805 the necessitation claim do not establish the necessitation claim. Those who  
1806 think that the necessitation claim is true, think that it is true on the evaluative  
1807 sense of “rational,” but for independent reasons.<sup>11</sup> I suspect that no argument  
1808 could establish the necessitation claim on the dispositional sense of “rational,”  
1809 so I am sceptical about the prospects of the argument from necessitation.  
1810 Nonetheless, as I will now argue, there is a better argument for constitutivism  
1811 about introspective belief available, one which is implicit in Shoemaker’s  
1812 work but which gets overshadowed by the argument from necessitation.

### 1813 **The Arguments from Cognitive Dynamics**

1814 In this section, I will formulate and evaluate the argument from cognitive  
1815 dynamics for constitutivism about instrumental desire and the argument from  
1816 cognitive dynamics for constitutivism about introspective belief. I will argue  
1817 that the former is a good argument for constitutivism about instrumental de-  
1818 sire, and, in light of the conclusions of the previous section, the *best* argument  
1819 for this thesis. Then I will argue that the latter is a good argument, if not an  
1820 equally good argument, for constitutivism about introspective belief.

#### 3.1 *For Constitutivism About Instrumental Desire*

1822 The argument from cognitive dynamics for constitutivism about instrumental  
1823 desire begins from an observation about the cognitive dynamics of instrumen-

---

10 The supervenience claim is the conclusion of Shoemaker’s famous argument from self-blindness. See Shoemaker (1996, 47–48).

11 Kieran Setiya argues that a nearby necessitation claim is true, and that this is what Shoemaker has correctly drawn attention to (Setiya 2011). But what gets necessitated, according to Setiya, is a capacity for introspective belief, not introspective belief itself.

1824 tal desires. The observation is that they systematically come and go with the  
1825 desires and means-end beliefs which explain them.

1826  $D_1$ . (i) If someone comes to desire to  $\Psi$  and to believe that  $\Phi$ -ing  
1827 is a means to  $\Psi$ -ing, they will come to desire to  $\Phi$  if they are fully  
1828 rational and (ii) if someone merely instrumentally desires to  $\Phi$  and  
1829 they cease desiring to  $\Psi$  or cease believing that  $\Phi$ -ing is a means to  
1830  $\Psi$ -ing, then they will cease desiring to  $\Phi$  if they are fully rational.

1831 These are claims about the cognitive dynamics of particular desires. They  
1832 are the cognitive dynamics of someone who is rational, that is, someone  
1833 who is disposed to conform to the principles of rationality. These cognitive  
1834 dynamics are partly constitutive of what it is to be disposed to conform to the  
1835 principles of rationality. It is because we observe these cognitive dynamics  
1836 that we believe that we are rational in this sense.

1837 How could a cognitive system exhibit these dynamics? To see how, let's  
1838 assume a broadly functionalist picture of beliefs and desires. On this picture,  
1839 for  $X$  to believe that  $P$  is for  $X$  to be in some state or other which plays the  
1840 believing-that- $P$ -role, and for  $X$  to desire to  $\Phi$  is for  $X$  to be in some state  
1841 or other which plays the desiring-to- $\Phi$ -role. We could then re-describe the  
1842 cognitive dynamics above in terms of ceasing to be in a state which plays  
1843 the desiring-to- $\Phi$ -role under certain conditions, and coming to be in a state  
1844 which plays the desiring-to- $\Phi$ -role under certain conditions. But now two  
1845 importantly different hypotheses arise concerning the relations between these  
1846 states.

1847 According to one hypothesis, call it the *causal hypothesis*, the dynam-  
1848 ics are explained by the fact that when someone comes to be in a state  
1849 which plays the desiring-to- $\Phi$ -role, and someone comes to be in a state which  
1850 plays the believing-that- $\Psi$ -ing-is-a-means-to- $\Phi$ -ing-role, their coming to be  
1851 in these states jointly *causes* them—by means of their manifesting a rational  
1852 disposition—to come to be in a state which plays the desiring-to- $\Psi$ -role. Simi-  
1853 larly, someone's ceasing to be in a state which plays the believing-that- $\Psi$ -ing-is-  
1854 a-means-to- $\Phi$ -ing-role or ceasing to be in a state which plays the desiring-to- $\Phi$ -  
1855 role will cause them to cease being in a state which plays the desiring-to- $\Psi$ -role.  
1856 Rational dispositions can then be thought of as ordinary causal dispositions,  
1857 where the triggering conditions are thought of as causes of the manifestations.  
1858 On this hypothesis the cognitive dynamics are explained by various causal  
1859 transactions between constitutively distinct states or events involving consti-

1860 tutively distinct states (they must be constitutively distinct in order to stand  
1861 in causal relations).

1862 According to another hypothesis, call it the *constitutive hypothesis*, the  
1863 dynamics are explained by the fact that when someone comes to be in a  
1864 state which plays the desiring-to- $\Phi$ -role and someone comes to be in a state  
1865 which plays the believing-that- $\Psi$ -ing-is-a-means-to- $\Phi$ -ing-role, their coming  
1866 to be in these states constitutes—by means of their manifesting a rational  
1867 disposition—their coming to be in a state which plays the desiring-to- $\Psi$ -role.  
1868 It does so because the former states together *constitute* a state which plays  
1869 the desiring-to- $\Phi$ -role. Similarly, someone's ceasing to be in a state which  
1870 plays the believing-that- $\Psi$ -ing-is-a-means-to- $\Phi$ -ing-role or ceasing to be in a  
1871 state which plays the desiring-to- $\Phi$ -role just is their ceasing to be in a state  
1872 which plays the desiring-to- $\Psi$ -role. Rational dispositions, on this hypothesis,  
1873 are not ordinary causal dispositions. They are what we might call constitutive  
1874 dispositions, since the triggering conditions bear a constitutive relation to the  
1875 manifestations. On this hypothesis the cognitive dynamics are explained in  
1876 terms of states which play particular roles jointly constituting states which  
1877 play other roles. The relevant thing about the constitutive hypothesis is that,  
1878 if it is true, then instrumental desires are not distinct from the desires and  
1879 means-end beliefs which explain them. This is because instrumental desires  
1880 are partly constituted by states which partly constitute the corresponding  
1881 desires and means-end beliefs which explain them. So if there is an argument  
1882 for the constitutive hypothesis, there is an argument for constitutivism about  
1883 instrumental desire.

1884 Each of these hypotheses is clearly an empirical hypothesis. If cognitive  
1885 science were so advanced that we could determine which states play which  
1886 roles, then, in principle we could settle the question of whether the states  
1887 are constitutively distinct and causally related or constitutively non-distinct  
1888 and constitutively rather than causally related. But we are far from being  
1889 able to answer the question this way. The best we have, and the best we may  
1890 ever have, is indirect evidence for one hypothesis over the other based on  
1891 arguments to the best explanation of the observed cognitive dynamics. Let's  
1892 consider the evidence for and against, then.

1893 Perhaps the weakest consideration in favour of the constitutive hypothesis,  
1894 but a consideration nonetheless, comes from the relative cognitive efficiency  
1895 of having the states which play the role of certain desires being constituted by  
1896 the states which play the roles of other desires and means-end beliefs, rather  
1897 than having the former be distinct from and caused by the latter. A cognitive

1898 system requires far fewer distinct states and fewer dependencies between  
1899 them in order to have a wide range of instrumental desires on the constitutive  
1900 hypothesis. To put the point in slogan form: *the constitutive hypothesis is*  
1901 *cognitively more efficient than the causal hypothesis.*

1902 A stronger consideration in favour of the constitutive hypothesis comes  
1903 from the observation that the causal roles which are definitive of desires  
1904 and means-end beliefs, along with the principles of rationality, predict that  
1905 by merely having those desires and means-end beliefs, and being rational,  
1906 the agent will be disposed to act *as if* she desired the means. If desiring the  
1907 means were a matter of coming to be in a distinct state which plays the role  
1908 of desiring the means, as the causal hypothesis holds, then the disposition to  
1909 act as if one desired the means would be over-determined. There's nothing by  
1910 way of the agent's dispositions to act that being in this state would contribute  
1911 which is not already contributed by their being in these other states and their  
1912 being rational. The state is motivationally redundant.<sup>12</sup> So, to put the point in  
1913 slogan form: *the constitutive hypothesis avoids the prediction that instrumental*  
1914 *desires are motivationally redundant.*

1915 Perhaps the strongest consideration in favour of the constitutive hypothesis,  
1916 the one that I am willing to put the most weight on, begins with an observation  
1917 about the strength of the dependence of instrumental desires on means-end  
1918 beliefs and other desires. As we saw earlier, one claim about the cognitive  
1919 dynamics of instrumental desires is that if someone merely instrumentally  
1920 desires to  $\Phi$  and they cease desiring to  $\Psi$  or cease believing that  $\Phi$ -ing is  
1921 a means to  $\Psi$ -ing, then they will cease desiring to  $\Phi$ . While this claim is  
1922 compatible with both the constitutive and causal hypotheses, the constitutive  
1923 hypothesis has a far better explanation of it. Indeed, the explanation comes for  
1924 free on the constitutive hypothesis, since it is no surprise that when one ceases  
1925 to be in either of the states which jointly constitutes the state which plays the  
1926 desiring-to- $\Phi$ -role that one will cease desiring to  $\Phi$ . The causal hypothesis  
1927 requires the auxiliary hypothesis here that when the state which plays the  
1928 desiring-to- $\Phi$ -role is caused and explained by the states which play the other  
1929 roles, it will remain causally dependent on those states. Of course, it is possible  
1930 that, simply as a matter of fact, such a state will remain causally dependent  
1931 on these other states. But this claim has to be added as an auxiliary hypothesis  
1932 to the causal hypothesis, thus making the hypothesis more complicated than  
1933 the constitutive hypothesis. And, moreover, there is every reason to think

---

12 See Arpaly and Schroeder (2014, 9) for a similar observation.

1934 that this causal dependence would sometimes break down, giving rise to  
 1935 stray instrumental desires, desires which are no longer dependent on the  
 1936 means-end beliefs and desires which caused them in the first place. But if this  
 1937 phenomenon exists in our psychology, it remains unobserved. The constitutive  
 1938 hypothesis correctly predicts that there will be no stray instrumental desires.  
 1939 To sum up the point of this paragraph in slogan form: *instrumental desires are*  
 1940 *deeply dependent on other mental states for their existence and the constitutive*  
 1941 *hypothesis best explains this.*

1942 It is basically this consideration which motivates the argument in the fol-  
 1943 lowing passage from Smith:

1944 It is a striking fact that instrumental desires disappear immedi-  
 1945 ately an agent loses either the relevant non-instrumental desire  
 1946 or means-end belief [...]. Yet there is no reason why this should  
 1947 be so if an instrumental desire were merely a desire that has a  
 1948 non-instrumental desire and a means-end belief somewhere in  
 1949 its causal history. Why should a desire disappear when (say) the  
 1950 desire that caused it, way back when, disappears? Instrumental  
 1951 desires are thus better thought of as being nothing over and above  
 1952 the non-instrumental desires and means-end beliefs that explain  
 1953 them. (2004, 96)

1954 Smith begins here with the observation that instrumental desires are deeply  
 1955 dependent on the desires and means-end beliefs which explain them. He then  
 1956 argues against the causal hypothesis and for the constitutive hypothesis on the  
 1957 basis of the fact that the latter provides a better explanation of the observation  
 1958 than the former. Of course, Smith's argument against the causal hypothesis is  
 1959 too fast. An instrumental desire could remain causally dependent on another  
 1960 desire and a means-end belief in the way that a light's being on remains  
 1961 causally dependent on the light switch's being turned on. It could be that  
 1962 the relation between an instrumental desire and the desire and means-end  
 1963 belief which explains it is like this. But the constitutive hypothesis nonetheless  
 1964 provides a better explanation of the deep dependence between an instrumental  
 1965 desire and the desire and means-end belief which explains it.

1966 The best argument, then, for constitutivism about instrumental desire is  
 1967 the argument from cognitive dynamics. Unlike the argument from necessi-  
 1968 tation for constitutivism about instrumental desire, which would decisively  
 1969 establish constitutivism about instrumental desire if it were cogent, the ar-  
 1970 gument from cognitive dynamics makes constitutivism about instrumental



1971 desire the conclusion of an ordinary argument to the best explanation. But  
 1972 it isn't really surprising that this should be so, since constitutivism about  
 1973 instrumental desire is most plausibly thought of as a contingent hypothesis  
 1974 about instrumental desires.

### 3.2 For Constitutivism About Introspective Belief

1976 If the argument from cognitive dynamics for constitutivism about instrumen-  
 1977 tal desire provides good support for constitutivism about instrumental desire,  
 1978 then perhaps an analogous argument from cognitive dynamics could provide  
 1979 good support for constitutivism about introspective belief. In this section I  
 1980 argue that it does.

1981 The argument from cognitive dynamics for constitutivism about intro-  
 1982 spective belief begins from an observation about the cognitive dynamics of  
 1983 introspective beliefs. The observation is that introspective beliefs systemati-  
 1984 cally come and go with the mental states—or in the case under consideration,  
 1985 the attitudes—that they are about.

1986  $D_2$ . (i) If someone has some interest in the question of whether  
 1987 they  $\Psi$  that P, and they come to  $\Psi$  that P<sup>13</sup>, then they will come  
 1988 to believe that they  $\Psi$  that P if they are fully rational, and (ii) if  
 1989 someone introspectively believes that they  $\Psi$  that P, and they cease  
 1990  $\Psi$ -ing that P, then they will cease believing that they  $\Psi$  that P if they  
 1991 are fully rational.

1992 These are claims about the cognitive dynamics of particular beliefs. They  
 1993 are the cognitive dynamics of someone who is rational, that is, someone  
 1994 who is disposed to conform to the principles of rationality. These cognitive  
 1995 dynamics are partly constitutive of what it is to be disposed to conform to the  
 1996 principles of rationality. It is because we observe these cognitive dynamics  
 1997 that we believe that we are rational in this sense.

1998 Now, let me be upfront here about an important disanalogy with the obser-  
 1999 vation about the cognitive dynamics of instrumental desires. The observation  
 2000 just given is likely to strike many as highly controversial. Failures of intro-  
 2001 spective belief are the norm. Failures of instrumental desire are the exception.  
 2002 I have two responses to this kind of pessimism about introspective belief. The

13 The antecedent of this conditional may need strengthening. Perhaps one needs to consider the question of whether one  $\Psi$ -s that P also.

2003 first is that it vastly overstates the case. To say that failures of introspective  
 2004 belief are the norm is to overlook the wide range of cases where introspective  
 2005 belief is utterly unproblematic. Over a vast range of mundane beliefs, desires,  
 2006 and other mental states, I have utterly unproblematic introspective access.  
 2007 If failures of introspective belief were the norm here, our mental life would  
 2008 be in serious trouble. The second response is that nothing in the argument  
 2009 from cognitive dynamics depends on an overly optimistic view of our capacity  
 2010 for introspective belief. It may well be that failures of local rationality are far  
 2011 more common in the case of introspective beliefs, but as long as introspective  
 2012 beliefs have some of the features I draw attention to below, the argument  
 2013 from cognitive dynamics will go through. This is a point I will return to after  
 2014 presenting the rest of the argument.

2015 How could a cognitive system exhibit the dynamics above? Again, assum-  
 2016 ing a broadly functionalist picture of beliefs and desires, we can recast the  
 2017 observation in the following terms. When someone comes to be in a state  
 2018 which plays the  $\Psi$ -ing-that-P-role, they come to be in a state which plays the  
 2019 believing-that-one- $\Psi$ -s-that-P-/-to- $\Phi$ -role, insofar as they are rational and have  
 2020 some interest in the question of whether they  $\Psi$  that P. When someone ceases  
 2021 to be in a state which plays the  $\Psi$ -ing-that-P-role, they cease to be in a state  
 2022 which plays the believing-that-one- $\Psi$ -s-that-P-role. Again, two importantly  
 2023 different hypotheses arise concerning the identity of these states.

2024 According to one hypothesis, call it the *causal hypothesis*, the dynamics are  
 2025 explained by the fact that when someone comes to be in a state which plays  
 2026 the  $\Psi$ -ing-that-P-role, they are caused to come to be in a state which plays the  
 2027 believing-that-one- $\Psi$ -s-that-P-role, insofar as they are rational and have some  
 2028 interest in the question of whether they  $\Psi$ -that-P. And someone ceases to be  
 2029 in a state which plays the  $\Psi$ -ing-that-P-role, they are caused to cease to be in  
 2030 a state which plays the believing-that-one- $\Psi$ -s-that-P-role. On this hypothesis  
 2031 the cognitive dynamics are explained by various causal transactions between  
 2032 constitutively distinct states or events involving these states.

2033 According to another hypothesis, call it the *constitutive hypothesis*, the  
 2034 dynamics are explained by the fact that when someone comes to be in a  
 2035 state which plays the  $\Psi$ -ing-that-P-role, they thereby come to be in a state  
 2036 which plays the believing-that-one- $\Psi$ -s-that-P-role, since the former state  
 2037 plays the latter role. And when someone ceases to be in a state which plays  
 2038 the  $\Psi$ -ing-that-P-role, they thereby cease being in a state which plays the  
 2039 believing-that-one- $\Psi$ -s-that-P-role, since it was the former state which played

2040 the latter role. On this hypothesis the cognitive dynamics are explained by  
2041 the states which play the first-order roles playing the second-order roles.<sup>14</sup>

2042 Each of these hypotheses is clearly an empirical hypothesis. What argu-  
2043 ments can be given for and against? Not surprisingly, the considerations are  
2044 perfectly analogous to those given in the argument from cognitive dynamics  
2045 for constitutivism about instrumental desire.

2046 Perhaps the weakest consideration in favour of the constitutive hypothesis,  
2047 but a consideration nonetheless, comes from the relative cognitive efficiency  
2048 of having the states which play the roles of particular mental states also  
2049 play the role of beliefs about those mental states, rather than having the  
2050 former be distinct from and caused by the latter. A cognitive system requires  
2051 far fewer distinct states and fewer dependencies between them in order to  
2052 have a wide range of introspective beliefs on the constitutive hypothesis. This  
2053 consideration has considerable bite in contemporary contexts where doubt has  
2054 arisen, both on the basis of philosophical and empirical enquiry, concerning  
2055 the claim that we have a distinct perception-like capacity for inner-sense, one  
2056 which causally detects our mental states and outputs introspective beliefs.  
2057 This view goes hand in hand with the causal hypothesis. The constitutive  
2058 hypothesis cuts out the middle-man, and requires no distinct perception-like  
2059 capacity for inner-sense.<sup>15</sup> To sum up the points of this paragraph in slogan  
2060 form: *the constitutive hypothesis is cognitively more efficient than the causal*  
2061 *hypothesis.*

2062 A stronger consideration in favour of the constitutive hypothesis comes  
2063 from the observation that the causal roles which are definitive of many mental  
2064 states, along with the principles of rationality, predict that by merely being  
2065 in those mental states, and being rational, an agent will be disposed to act  
2066 *as if* she believes that she is in those mental states.<sup>16</sup> In particular, it has  
2067 been observed that if you are in pain, say, and you are rational, you will be  
2068 disposed, partly in virtue of the fact that you are in pain, to say “I am in pain.”

---

14 See Shoemaker (1996, 33–34, 242–244).

15 This might be a little bit unfair to the causal hypothesis. Since there is a version of the causal hypothesis which cuts out the middle-man too, and requires no distinct perception-like capacity for inner-sense. My point here is that once we are on the lookout for cognitively efficient hypotheses about introspective beliefs, the constitutive hypothesis wins hands down.

16 This point is well made by Shoemaker (1996). While Shoemaker makes the point in the context of arguing for the necessitation or supervenience claim, I am here making it in the context of the cognitive dynamics argument. Shoemaker sometimes says things which suggest that he might have something like the cognitive dynamics argument in mind. This is a point I will come to in the text.

2069 This point has been made over and over by expressivists in the philosophy  
 2070 of introspection, who take it to show that statements like “I am in pain” do  
 2071 not report mental states but merely express them. But we needn’t understand  
 2072 the claim in this manner, since it is possible that such statements both report  
 2073 and express mental states, because the mental states they express double,  
 2074 according to the constitutive hypothesis, as beliefs about those mental states.  
 2075 If introspectively believing that you are in a particular mental state were  
 2076 a matter of coming to be in a distinct state which plays the role of a belief  
 2077 that you are in some mental state, as the causal hypothesis holds, then the  
 2078 disposition to act as if you believed that you were in that mental state would  
 2079 be over-determined. There’s nothing by way of the agent’s dispositions to  
 2080 act that being in this distinct state would contribute which is not already  
 2081 contributed by their being in the first-order state, their having an interest in  
 2082 the question of whether they are in that state, and their being fully rational.  
 2083 The state is motivationally redundant.<sup>17</sup> So, to put the point in slogan form:  
 2084 *the constitutive hypothesis avoids the prediction that introspective beliefs are*  
 2085 *motivationally redundant.*

2086 Perhaps the strongest consideration in favour of the constitutive hypothesis,  
 2087 the one that I am willing to put the most weight on, begins with an obser-  
 2088 vation about the strength of the dependence of introspective beliefs on the  
 2089 mental states they are about. As we saw earlier, one claim about the cognitive  
 2090 dynamics of introspective belief is that if someone introspectively believes  
 2091 that they  $\Psi$  that P, and they cease  $\Psi$ -ing that P, then they will cease believing  
 2092 that they  $\Psi$  that P.<sup>18</sup> While this claim is compatible with both the constitutive  
 2093 and causal hypotheses, the constitutive hypothesis has a far better explanation  
 2094 of it. Indeed, the explanation comes for free on the constitutive hypothesis,  
 2095 since it is no surprise that when one ceases to be in the state which plays the  
 2096 believing-that-one- $\Psi$ -s-that-P-role, one will cease believing that one  $\Psi$ s that  
 2097 P. The causal hypothesis requires the auxiliary hypothesis here that when  
 2098 the state which plays the believing-that-one- $\Psi$ -s-that-P-role is caused and  
 2099 explained by the state which plays the  $\Psi$ -that-P-role, it will remain causally

---

17 To be clear: I am not denying that introspective beliefs themselves are motivationally redundant, that they make no difference to the cognitive functioning of the mental states they are about; I am only claiming that a distinct state which had these consequences would be redundant insofar as being in the mental states in question, having an interest in the question of whether you are in those mental states, and being fully rational would already have the consequences for one’s cognitive life we take introspective beliefs to have.

18 This was pointed out to me by Daniel Nolan.

2100 dependent on this state. Of course, it is possible that, simply as a matter of  
2101 fact, such a state will remain causally dependent on this other state. But this  
2102 claim has to be added as an auxiliary hypothesis to the causal hypothesis, thus  
2103 making the hypothesis more complicated than the constitutive hypothesis.  
2104 And, moreover, there is every reason to think that this causal dependence  
2105 would sometimes break down, giving rise to stray introspective beliefs, *intro-*  
2106 *spective* beliefs which are no longer dependent on the mental states which  
2107 caused them in the first place. But if this phenomenon exists in our psychol-  
2108 ogy, it remains unobserved. The constitutive hypothesis correctly predicts that  
2109 there will be no stray introspective beliefs. And this is a significant point in its  
2110 favour. So, to sum up the points of this paragraph in slogan form: *introspective*  
2111 *beliefs are deeply dependent on other mental states for their existence and the*  
2112 *constitutive hypothesis best explains this.*

2113 While Shoemaker is more closely associated with the argument from nec-  
2114 cessitation, it is clear that he also has something like the argument from  
2115 cognitive dynamics in mind. Indeed, I think that, to the extent that he does  
2116 have the latter in mind, this is the best argument he has for constitutivism  
2117 about introspective belief. Consider the following passage from a recent paper  
2118 of Shoemaker's defending constitutivism about introspective belief:

2119 One might, indeed, wonder whether there is any need to postulate  
2120 standing second-order beliefs that self-ascribe available first-order  
2121 beliefs. It goes with having the available first-order belief that *p*  
2122 that if the question whether one believes that *p* arises, one will  
2123 judge that one does—one will assent to the proposition that one  
2124 believes that *p*. But this seems to be the result of one's having the  
2125 belief that *p*, not the result of one's having a second-order belief  
2126 whose cognitive dynamics is independent of that of the belief that  
2127 *p*, in the way that the cognitive dynamics of one's belief about  
2128 another person's belief is independent of that of the other person's  
2129 belief. It would seem inefficient for our psychology to involve the  
2130 storage of standing second-order beliefs ascribing available first-  
2131 order beliefs, if there is nothing for these second-order beliefs to  
2132 do that is not done by the first-order beliefs themselves. (2012,  
2133 247)

2134 These remarks combine elements of all three of the considerations I have  
2135 given above. Shoemaker makes a claim about efficiency, there is also a claim

2136 about redundancy, and he speaks of the predicted cognitive independence  
 2137 of introspective beliefs from the mental states they are about on the causal  
 2138 hypothesis. I have teased out these considerations and argued that together  
 2139 they add up to a reasonable case for the constitutive hypothesis.

2140 The best argument, then, for constitutivism about introspective belief is  
 2141 the argument from cognitive dynamics. Unlike the argument from necessi-  
 2142 tation for constitutivism about introspective belief, which would decisively  
 2143 establish constitutivism about introspective belief if it were cogent, the ar-  
 2144 gument from cognitive dynamics makes constitutivism about introspective  
 2145 belief the conclusion of an ordinary argument to the best explanation. But  
 2146 it isn't really surprising that this should be so, since constitutivism about  
 2147 introspective belief is most plausibly thought of as a contingent hypothesis  
 2148 about introspective beliefs.

2149 To end, let me return to the obvious line of criticism which may be raised  
 2150 against the argument from cognitive dynamics for constitutivism about in-  
 2151 trospective belief. The criticism is basically that it depends on far too rosy a  
 2152 picture of introspective belief. But we are now in a position to see that it does  
 2153 not. The consideration about efficiency requires only that we have a signifi-  
 2154 cant number of introspective beliefs, so that considerations of efficiency come  
 2155 into play. It doesn't require that we approximate omniscience and infallibility.  
 2156 The considerations about redundancy, likewise, only require that we have  
 2157 a significant number of introspective beliefs, so that considerations about  
 2158 redundancy come into play. And, finally, considerations about deep depen-  
 2159 dence do not require that we are rarely in error about our own mental states.  
 2160 This is an important point. We may have many false beliefs about our own  
 2161 mental states. We may often be in error about our own mental states. But as  
 2162 long as those beliefs about our own mental states arrived at by introspection  
 2163 remain deeply dependent on the mental states they are about—that is, as long  
 2164 as there are no stray *introspective* beliefs—the point about dependence holds.\*

2165

---

\* Thanks to Michael Smith for discussion on a very early version of this paper. Thanks also to members of the ANU Philosophy of Mind Work in Progress Group for their helpful feedback on an earlier version of this paper. Finally, thanks to several reviewers for *Dialectica* for their helpful comments.

Ryan Cox

 0000-0002-1381-448X

The University of Sydney

ryan.cox@sydney.edu.au

2166  
2167  
2168  
21692170 **References**

- 2171 ARMSTRONG, David M. 1968. *A Materialist Theory of the Mind*. London: Routledge &  
2172 Kegan Paul. Reissued as Armstrong (2023).
- 2173 —. 2023. *A Materialist Theory of the Mind*. London: Routledge. With a new foreword  
2174 by Peter Anstey, doi:10.4324/b23154.
- 2175 ARPALY, Nomy and SCHROEDER, Timothy. 2014. In *Praise of Desire*. Oxford: Oxford  
2176 University Press, doi:10.1093/acprof:oso/9780199348169.001.0001.
- 2177 BOYLE, Matthew. 2009. “Two Kinds of Self-Knowledge.” *Philosophy and Phenomeno-*  
2178 *logical Research* 78(1): 133–164, doi:10.1111/j.1933-1592.2008.00235.x.
- 2179 BYRNE, Alex. 2011. “Transparency, Belief, Intention.” *Proceedings of the Aristotelian*  
2180 *Society, Supplementary Volume* 85: 201–221, doi:10.1111/j.1467-8349.2011.00203.x.
- 2181 CASSAM, Quassim. 2014. *Self-Knowledge for Humans*. Oxford: Oxford University Press,  
2182 doi:10.1093/acprof:oso/9780199657575.001.0001.
- 2183 CHALMERS, David J. 2010. *The Character of Consciousness*. Oxford: Oxford University  
2184 Press, doi:10.1093/acprof:oso/9780195311105.001.0001.
- 2185 DAVIS, Wayne A. 1986. “The Two Senses of Desire.” in *The Ways of Desire: New Essays*  
2186 *in Philosophical Psychology on the Concept of Wanting*, edited by Joel MARKS, pp.  
2187 63–82. Chicago, Illinois: Precedent Publishing Inc.
- 2188 GERTLER, Brie. 2010. *Self-Knowledge*. London: Routledge, doi:10.4324/9780203835678.
- 2189 GORDON, Robert M. 1986. “Folk Psychology as Simulation.” *Mind and Language* 1(2):  
2190 158–171, doi:10.1111/j.1468-0017.1986.tb00324.x.
- 2191 HUME, David. 1748. *Philosophical Essays Concerning Human Understanding*. London:  
2192 Andrew Millar of the Strand.
- 2193 —. 1975. *Enquiries Concerning Human Understanding and Concerning the Principles*  
2194 *of Morals*. 3rd ed. Oxford: Oxford University Press. Edited by Lewis Amherst Selby-  
2195 Bigge; revised and with notes by Peter Harold Nidditch, doi:10.1093/oseo/instanc  
2196 e.00046349.
- 2197 —. 1999. *An Enquiry Concerning Human Understanding*. Oxford: Oxford University  
2198 Press. Edited by Tom L. Beauchamp, doi:10.1093/owc/9780199549900.001.0001.
- 2199 KIND, Amy. 2003. “Shoemaker, Self-Blindness and Moore’s Paradox.” *The Philosophical*  
2200 *Quarterly* 53(210): 39–48, doi:10.1111/1467-9213.00294.
- 2201 MARKS, Joel. 1986. “Introduction: On the Need for a Theory of Desire.” in *The Ways of*  
2202 *Desire: New Essays in Philosophical Psychology on the Concept of Wanting*, edited  
2203 by Joel MARKS, pp. 1–15. Chicago, Illinois: Precedent Publishing Inc.

- 2204 MCDANIEL, Kris and BRADLEY, Ben. 2008. "Desires." *Mind* 117(466): 267–302, doi:10  
 2205 .1093/mind/fzn044.
- 2206 MORAN, Richard. 2001. *Authority and Estrangement: An Essay on Self-Knowledge*.  
 2207 Princeton, New Jersey: Princeton University Press, doi:10.1515/9781400842971.
- 2208 SCHROEDER, Timothy. 2004. "Functions from Regulation." *The Monist* 87(1): 115–135,  
 2209 doi:10.5840/monist20048717.
- 2210 SETIYA, Kieran. 2011. "Knowledge of Intention." in *Essays on Anscombe's Intention*,  
 2211 edited by Anton FORD, Jennifer HORNSBY, and Frederick STOUTLAND, pp. 170–197.  
 2212 Cambridge, Massachusetts: Harvard University Press, doi:10.4159/harvard.9780  
 2213 674060913.c7.
- 2214 SHOEMAKER, Sydney S. 1996. *The First-Person Perspective and Other Essays*. Cambridge:  
 2215 Cambridge University Press, doi:10.1017/cbo9780511624674.
- 2216 —. 2012. "Self-Intimation and Second-Order Belief." in *Introspection and Conscious-*  
 2217 *ness*, edited by Declan SMITHIES and Daniel STOLJAR, pp. 239–258. Oxford: Oxford  
 2218 University Press, doi:10.1093/acprof:oso/9780199744794.003.0009.
- 2219 SMITH, Michael A. 2004. "Instrumental Desires, Instrumental Rationality." *Proceedings*  
 2220 *of the Aristotelian Society, Supplementary Volume* 78: 93–109, doi:10.1111/j.0309-  
 2221 7013.2004.00117.x.
- 2222 STOLJAR, Daniel. 2019. "Evans on Transparency: A Rationalist Account." *Philosophical*  
 2223 *Studies* 176(8): 2067–2085, doi:10.1007/s11098-018-1111-x.



# The Mental States First Theory of Promising

ALIDA LIBERMAN

Most theories of promising are insufficiently broad, for they ground promissory obligation in some external or contingent feature of the promise. In this paper, I introduce a new kind of theory. The Mental States First (MSF) theory grounds promissory obligation in something internal and essential: the mental state expressed by promising, or the state that promisors purport to be in. My defense of MSF relies on three claims. First, promising to  $\Phi$  expresses that you have resolved to  $\Phi$ . Second, resolving to  $\Phi$  commits you to  $\Phi$ ing, all else being equal. Third, the norms on speech acts are determined by the norms on the mental states they express, such that publicly expressing that you are in a state subjects you to whatever commitments are normally incurred by being in that state, regardless of whether you really are in it. I suggest that this general approach might also explain how the norms on other sorts of speech acts work.

Philosophers have offered a variety of theories of promissory obligation, most of which ground promissory obligation in some external or contingent feature of the promise. In this paper, I sketch a new kind of theory, which instead grounds promissory obligation in something internal and essential to promise-making. This *Mental States First* theory of promising (or MSF) posits that the norms on promises are fixed or determined by the norms on the promises' underlying mental states; the mental states are "first" in the sense that they are explanatorily prior. My aim is to show that MSF should be taken seriously, because it is grounded in plausible assumptions, can accommodate a wider range of cases than can most theories, and fruitfully situates promising as part of a broader pattern of speech acts that behave in similar ways.

I motivate the search for a new theory of promissory obligation in section 1 by pointing out that most theories face cases of apparently genuine promissory obligation that they cannot accommodate. In section 2, I lay out how MSF

2254 works. In sections 3, 4, 5, I sketch arguments for several claims on which the  
 2255 success of MSF depends; if these claims hold, then MSF is a viable theory of  
 2256 promissory obligation. My argument in section 5 highlights how the general  
 2257 MSF pattern holds for other speech act/mental state pairs, which suggests  
 2258 that a Mental States First approach might be productive for understanding  
 2259 the norms on other kinds of speech acts, as well.

## 2260 **1 Motivating MSF with Marginal Cases**

2261 Most theories of promissory obligation have no trouble explaining why we  
 2262 are obligated to keep our promises in paradigmatic cases, or those in which  
 2263 the promise is something the promisee wants and expects will occur, the  
 2264 promisor intends to bind herself to act as promised, harm would occur were  
 2265 the promise broken, and the promise occurs within some obvious social  
 2266 convention of promising. However, these theories struggle with explaining  
 2267 why we are obligated to keep promises in non-paradigmatic cases, which  
 2268 I call *marginal cases*.<sup>1</sup> These are cases in which there are good theoretical  
 2269 arguments—or a strong intuitive presumption with no good arguments to  
 2270 the contrary—that the case generates a genuine promissory obligation, but in  
 2271 which it is doubtful whether the theory can explain or accommodate the case  
 2272 as an instance of promissory obligation. Proponents of particular theories tend  
 2273 to wield such cases as dialectical weapons against each other (e.g., by arguing  
 2274 that view A is inadequate because it cannot capture case X, or that view B  
 2275 should be preferred because it can).<sup>2</sup> This approach presumes that theories  
 2276 should be able to explain every plausible case of promissory obligation, which  
 2277 makes sense as a desideratum for a good account; more theoretical breadth  
 2278 and explanatory power is usually better than less. Moreover, since promises

---

1 My use of “marginal” to describe these cases is not meant to imply that these cases only barely count as promises. Rather, I mean to convey that they are cases that ought to be considered genuine promises but that lack some of the core features of the most obvious and paradigmatic promises, which places them at the edges or borders of our practices of promising. Thanks to an anonymous referee for discussion of this point.

2 For example, Scanlon criticizes conventionalist views in this way, by arguing that they cannot accommodate a proposed case of promising without a social practice. And Scanlon’s critics make similar sorts of arguments about his view, claiming that it fails because it is subject to a counterexample; see Cholbi (2002) and Southwood and Friedrich (2009). Of course, philosophers who argue in such ways typically also offer theoretical objections against the view they are targeting. But they generally assume that being unable to accommodate a case is highly problematic for a view.

2279 appear to be a unified phenomenon, we want a unified account of them if  
 2280 possible.

2281 It is important to note that a promise does *not* count as marginal for a theory  
 2282 simply because the theory implies that there are circumstances under which  
 2283 breaking the promise is morally permissible, all-things-considered. For it is  
 2284 widely accepted that promises yield pro tanto moral obligations, which can be  
 2285 overridden if excusing conditions arise (such as needing to break the promise  
 2286 in order to satisfy a more important conflicting obligation). When I claim  
 2287 that a case is marginal for a theory, I am claiming that the theory struggles to  
 2288 explain why the promise in the case has *any* moral force, or why there is even  
 2289 a pro tanto obligation to keep it.

2290 In the rest of this section, I give a brief overview of how the most popu-  
 2291 lar theories of promissory obligation are subject to marginal cases.<sup>3</sup> We can  
 2292 sort views with contemporary traction into three broad categories.<sup>4</sup> First are  
 2293 conventionalist accounts like that of David Hume—according to which we  
 2294 have reason to keep promises because it is bad for us (and our reputations  
 2295 as trustworthy people) if we do not—or John Rawls, who argues that failing  
 2296 to keep promises problematically free-rides on a valuable social practice.<sup>5</sup>  
 2297 Second are expectationalist accounts like that of T.M. Scanlon, according to  
 2298 which promise-breaking is impermissible because it involves violating expect-  
 2299 ations that you have raised in the promisee.<sup>6</sup> A third class of recently popular

3 This is not meant to be an exhaustive survey of theories of promissory obligation. I do not address views without much popular support or influence in the contemporary literature on promissory obligation, such as virtue-ethical approaches, intuitionist accounts, and Kantian deontology. Nor do I discuss act consequentialist views, according to which promises do not generate pro tanto obligations, and promise-keeping is morally required only if this leads to the best overall consequences. Of necessity, my discussions of the theories that I do address are sketchy and superficial; for more detail on how these (and other kinds of) theories of promissory obligation are subject to marginal cases, see Liberman (2015, chap. 2).

4 This follows Heuer's (2012) helpful classification.

5 See Hume (T 3.2.5.10) and Rawls (1971). More recent conventionalist accounts vary greatly in their details; for example, Sheinman (2008) argues that promisors can give themselves social practice-based reasons to act by communicating the intention to give themselves such a reason. Hooker (2011) offers a rule consequentialist view that can be understood as a conventionalist view, which grounds promissory obligation in the rules that would lead to the best consequences were people to internalize them (see Liberman 2020 for further discussion of this view). A number of theorists have also proposed hybrid conventionalist/expectationalist views, which ground promissory obligation in expectations that can only be generated from within an existing social practice; see Kolodny and Wallace (2003).

6 Scanlon invokes a principle of fidelity that is justifiable on contractualist grounds; see (1990) and chapter 7 of (1998). Other views of the same sort ground promissory obligation in something

2300 views are normative power accounts, which posit that we have the ability to  
 2301 change the normative situation directly and by declaration by exercising a  
 2302 normative power, through which we “change what someone is obliged to do  
 2303 by intentionally communicating the intention of hereby so doing” (Owens  
 2304 2012, 4). Normative power theories claim that promissory obligation stems  
 2305 from an exercise of such power, and that these powers exist because they are  
 2306 valuable for us to have. In general, these views claim that X is an important  
 2307 feature of our normative lives, and that we can possess X only by having the  
 2308 ability to bind ourselves to each other through promising. In a transcendental  
 2309 step, they conclude that we therefore must possess the ability to bind ourselves  
 2310 to each other through promising, lest X be inaccessible.<sup>7</sup>

2311 Each of these types of view—conventionalist, expectationalist, and norma-  
 2312 tive power—gets something right, as each captures an important feature of  
 2313 paradigmatic promise-making, a feature which tells us something about why  
 2314 we are obligated to keep promises in many or even most cases. But each of  
 2315 these theories is also in some way incomplete; while they can explain why a  
 2316 moral obligation is present in paradigmatic cases of promising, they cannot ex-  
 2317 plain the full breadth of cases in which we are obligated to keep our promises.  
 2318 For example, promise-breaking is often against one’s self-interest because of  
 2319 the moral sanction it incurs, as Hume suggests. But we can easily imagine  
 2320 cases where this is not so; consider a traveler passing through a remote town  
 2321 who cons locals out of their cash, promises to pay them back, and disappears.  
 2322 Similarly, most societies have robust promising conventions that it would be  
 2323 unfair to free-ride upon, as Rawls claims. But we can imagine a successful

---

other than expectations; for example, Judith Jarvis Thomson (1990) argues that promise-keeping is required because the promisee is relying on the promisor, while Daniel Friedrich and Nicholas Southwood (2009) argue that promise-breaking violates the trust you have invited the promisee to place in you. Cholbi (2002) offers a contractualist view that grounds promissory obligation not in what the promisee *actually* expects, but in what she is *entitled to* expect; I argue in Liberman (2015, chap. 2) that this view cannot accommodate cases in which the promisee lacks moral standing to hold the promisor accountable.

7 Normative power theorists cash out what the normative power to promise consists in and what valuable feature of our normative lives it supports in a variety of ways. For example, David Owens proposes a normative power grounded in what he calls our “authority interest,” or the interest we have in having a certain kind of practical authority over others. By making a promise, the promisor “give[s] the promisee the right to require performance of the promisor” (2012, 144); this serves our authority interest in allowing us to determine whether another person is obligated to act as she has promised, and in allowing others to make such determinations for us. Other normative power theorists speak instead of the ability to transfer rights (Shiffrin 2008) or create exclusionary reasons (Raz 1977).

2324 promise in a state of nature without such conventions; Scanlon offers a case of  
 2325 strangers from different societies on opposite sides of a river, whose hunting  
 2326 weapons have fallen on each other's sides and who manage to successfully  
 2327 promise to exchange the weapons in the absence of a shared convention of  
 2328 promising. Conventionalist views cannot accommodate such cases.<sup>8</sup>

2329 Similarly, Scanlon is right to point out that expectations matter; we are  
 2330 morally obligated not to mislead people and upset their expectations, or allow  
 2331 them to detrimentally rely on us and fail to follow through. But sometimes  
 2332 the promisee does not expect the promisor to perform; maybe the promisor is  
 2333 notoriously bad at keeping his word. As Berislav Marušić notes, expectational-  
 2334 ist views entail that if the promisor does “not succeed in forming expectations  
 2335 in the promisee [...] she will thus fail to incur promissory obligations; her  
 2336 promise, if it is one at all, won't be binding” (2013, 305). But making a promise  
 2337 that doesn't generate expectations shouldn't get the promisor off the moral  
 2338 hook; it would be problematic indeed if you were justified in promise-breaking  
 2339 because the promisee assumed you were unreliable, and accordingly did not  
 2340 expect you to keep your promise.

2341 Normative power accounts face their own marginal cases. For example,  
 2342 many people assume that there is a distinctively promissory obligation to  
 2343 keep promises we make to people who die before the promise can be kept.  
 2344 Promises to the dead are marginal for authority transfer views like Owens's.  
 2345 For there must be a person to whom the authority is transferred and in whom  
 2346 the authority continues to reside—but once the promisee is dead, no such  
 2347 authority bearer exists. Normative power theorists also presume that sincerely  
 2348 promising to  $\Phi$  requires intending to obligate yourself to  $\Phi$ . But this constraint  
 2349 leaves out some cases, as there seem to be cases of sincere promises that do not  
 2350 involve forming such an intention. As Thomas Pink (2009) argues, promisees  
 2351 are generally concerned not with whether the promisor is morally obligated  
 2352 to perform the promised action, but with whether the promisor *actually will*  
 2353 perform the promised action. And we frequently make promises without  
 2354 specifically intending to obligate ourselves in daily life.<sup>9</sup> I take it that you

8 More worryingly, Rawls's view fails to account for the directed nature of promissory obligation, as free-riding on a valuable social practice wrongs *all* of the participants in the practice equally, rather than wronging the promisee in particular. As Scanlon (1998) and Kolodny and Wallace (2003) both note, this is a major failure—although it is not a failure to accommodate a marginal case.

9 This is admittedly an overly quick dismissal of a popular view. I give a more detailed argument for this claim in Liberman (2015, chap. 4). See also Liberman (2015, chap. 5) for more about the sincerity conditions on promise-making.

2355 can promise to  $\Phi$  while merely *foreseeing* that your promise will incur an  
 2356 obligation but without specifically *intending* to incur such an obligation. For  
 2357 example, I can promise to take good care of a baseball glove I borrow from  
 2358 my brother with the intention of putting his mind at ease about getting the  
 2359 glove back in good condition, while merely *recognizing* (but not intending)  
 2360 that this promise will obligate me. I presume that you might also sincerely  
 2361 promise to  $\Phi$  without even foreseeing that you will  $\Phi$  (e.g., if you make a  
 2362 conditional promise and you believe that the relevant condition is extremely  
 2363 unlikely to occur), or to sincerely promise to  $\Phi$  while being a nihilist about  
 2364 whether obligations can ever be incurred. Normative power views cannot  
 2365 easily accommodate such cases.<sup>10</sup>

2366 The three kinds of theory of promissory obligation fail to explain why we are  
 2367 obligated to keep our promises in all—and only—cases of seemingly genuine  
 2368 promissory obligation. Conventionalist and expectationalist accounts risk  
 2369 being both over- and under-inclusive, because the sources of obligation that  
 2370 they point to are not *distinctively* promissory: that is, they stem not from the  
 2371 content or nature of the promise itself, but from some other, contingent feature  
 2372 of the situation (such as the context within which the promise occurs, or the  
 2373 downstream effects of the promise, which may not always be as expected).  
 2374 While promises often rely on social conventions, they need not do so, as  
 2375 we saw with Scanlon's weapon-exchange example. We could also have social  
 2376 conventions (such as an honor code) that led people to scrupulously keep their  
 2377 word for fear of social repercussion or free-riding without invoking promises  
 2378 at all. Likewise, not all genuinely binding promises raise expectations. And we  
 2379 sometimes raise others' expectations without making promises to them (e.g.,  
 2380 by putting on your coat and raising my expectation that you are about to leave  
 2381 the building). Normative power accounts purport to be grounded in a feature  
 2382 of promises as-such (i.e., the intention to obligate oneself that is expressed by  
 2383 a promise), but the feature they identify is not essential. Although forming

---

10 Normative power theorists generally construe exercises of normative power as intentions to change the normative situation (e.g., through promising, consenting, giving, etc.). It might be possible to develop a normative power view according to which the interpersonal act of promise-making itself counts as a direct exercise of normative power, regardless of what intentions (e.g., to obligate yourself, or simply to act) accompany the act of promise-making. Such a view could get around the specific marginal case I am currently addressing. But it would require positing a brute normative power without an underlying explanation of how this normative power functions, which would not be as deeply explanatory or satisfying an account of how and why promises bind than MSF can offer. Thanks to an anonymous referee for discussion of this point.

2384 an intention to obligate yourself and thereby transfer authority might be one  
2385 way of creating a promissory obligation, it is not necessary.

2386 This is not to say that theories of promissory obligation that are subject to  
2387 marginal cases are without worth. To the contrary, maintaining valuable social  
2388 practices and satisfying expectations matter morally. Considerations such as  
2389 these often make it the case that a particular promise is morally obligatory to  
2390 keep on independent grounds; people are often obligated to act as they have  
2391 promised for reasons having to do with fairness, harm avoidance, and the  
2392 like. But as already noted, such sources of obligation are neither maximally  
2393 broad nor distinctively promissory. In the rest of this paper, I argue that there  
2394 is another source of promissory obligation that is both broad enough to cover  
2395 all cases, and distinctive of promising as such.

## 2396 **2 Introducing the Mental States First Theory of Promising**

2397 The existence of marginal cases motivates grounding a theory of promissory  
2398 obligation not in the contingent and external circumstances of particular  
2399 promises, but in something that all promises share, something that is internal  
2400 and essential to the act of promising itself. One thing all promises have in  
2401 common is that they involve a communicative act, usually a verbal speech  
2402 act but sometimes a written or non-verbal act (for a nod or a stern look can  
2403 successfully communicate a promise in the right context).<sup>11</sup> Whenever you  
2404 perform such a communicative act, you *express* that you are in a certain mental  
2405 state, by which I mean that you convey to your audience that you are in that  
2406 state, and really are in it if your performance is sincere. If we can ground  
2407 promissory obligation in the expression of a mental state that all promises  
2408 share, we will have a maximally broad view.

2409 This brings us to the following theory:

2410 **MENTAL STATES FIRST THEORY OF PROMISING (MSF).** The obli-  
2411 gation to keep promises is derived from the norms on resolutions,  
2412 which are the mental state expressed by promising.

---

11 Could you have promises in which no direct communication occurs? In such cases of implicit promises, there is either successful indirect communication (e.g., if Sam and Chris agree to “go steady” in a context in which this conveys a commitment to monogamy, Chris breaks a promise by dating someone else), or there is no successful communication and no promise is made (e.g., if Sam and Chris start casually dating without any mutual presumption of sexual exclusivity, Sam does not break any implicit promise by dating someone else).

2413 MSF claims that one who promises to  $\Phi$  expresses that she resolves to  $\Phi$ ,  
 2414 and *in virtue of publicly expressing that resolution* is subsequently obligated  
 2415 to  $\Phi$ , all else being equal.<sup>12</sup> I am proposing that the norms on promising  
 2416 can be explained by appeal to the norms on the mental state expressed by  
 2417 promising. This is a claim about explanatory priority; the norm on the mental  
 2418 state explains the existence of the norm on the speech act, because it fixes or  
 2419 determines what that norm is.

2420 My argument for MSF is as follows:

- 2421 1. Promising to  $\Phi$  expresses that you have formed a resolution to  $\Phi$ , condi-  
 2422 tional on the promisee's acceptance.

---

12 Downie (1985) articulates an account of promising that shares some structural similarities with MSF. Downie argues that “a promise is essentially a matter of pledging oneself” and that the promisee’s “reliance and expectations are well-founded to the extent that they see that I already regard myself as obliged and they know me to be a man of my word [...]. To promise is always to state an intention in obligation-creating circumstances” (1985, 266). Downie argues that these circumstances are those in which “the intended projects have been made central and essential in one’s total concerns. The self has been identified with the projects, and carrying them out has become not only a moral obligation of practical consistency but a strong moral obligation of honour or self-fidelity” (1985, 269). That is, Downie argues that there is a moral duty to maintain self-consistency about projects with which you have identified your will, and the promisor has “identified his will with the project” and “pinned his self on the future as described in his pledged actions” in a way that makes “keeping of the promise essential to the preservation of his personal integrity,” such that “he will be diminished as a person if he breaks his word” (1985, 270). My view differs from Downie’s in two primary ways. First, Downie’s account accommodates only those promises that are deeply tied up with the promisor’s sense of self, while MSF explains promissory obligation more broadly, including in cases that are not very important and that your sense of integrity is not wrapped up in. Second, Downie grounds the moral force of promissory obligation in the need to adhere to your important resolutions. While I agree that promissory resolutions generate moral obligations, I argue in Section 4.2 that this is because of the way in which they are conditional on the acceptance of the promisee, and not because there are general moral duties of self-consistency. Michael Robins (1984) offers an account that draws on action theory to ground promissory obligation in intention—specifically, in the intention that the assent of the promisee will obligate you to act as you’ve promised to. Robins begins with an “irreducibly normative” (1984, 12) notion of intention that binds the intending agent to act in certain ways in the future. He argues that vows are intentions about which the agent cannot change their mind, and that promises are the transferal of this “exclusionary mandate” about how one will act to the “normative control of another person” (1984, 120). Robins argues that this transfer of the exclusionary mandate to the promisee transforms the requirement to abide by one’s vow into a moral obligation; various critics of Robins have argued that it is not clear what exactly this transfer consists in, how it occurs, or how it generates a moral obligation (see Cottingham 1985; Lemos 1987; Smith 1987). By contrast, MSF does not rely on any transfer (of an exclusionary mandate or any other right) to the promisee.



- 2423 2. Resolving to  $\Phi$  conditional on another person's acceptance rationally  
2424 and morally obligates you to  $\Phi$  (all else being equal).
- 2425 3. In general, publicly expressing that you are in state X obligates you to  
2426 act as X demands of you, regardless of whether you really are in state X.
- 2427 4. Therefore, promising to  $\Phi$  rationally and morally obligates you to  $\Phi$  (all  
2428 else being equal), regardless of whether you have really resolved to  $\Phi$ .

2429 The above argument depends on the following claims: (1) that promises  
2430 express resolutions; (2) that resolutions rationally (and in some cases, morally)  
2431 obligate you to act, all else being equal; and (3) the *determination claim*,  
2432 according to which the norms governing the mental state expressed by a  
2433 speech act (at least partially) determine what the norms on that speech act  
2434 are, in both sincere and insincere cases. The latter premises are not derived  
2435 from the former; these are all independent claims which together establish  
2436 the conclusion. In the next three sections, I offer arguments to support each  
2437 of these claims. Because I lack space to fully defend them at present, my  
2438 argument is conditional: *if* these claims are true, then MSF is an appealing  
2439 and viable theory of promissory obligation.

2440 To clarify, MSF states that the mental state expressed by a promise—i.e., a  
2441 resolution—is a necessary and ineliminable part of the explanation of why  
2442 we are morally required to keep our promises. The mental state alone is not  
2443 sufficient; an unarticulated mental state cannot generate interpersonal norms,  
2444 and so the public expression of that state plays an essential role, as well. But  
2445 the norms on the mental state are the original source of the obligation, and  
2446 they determine what the norm on the publicly expressed mental state is: the  
2447 norms on the promise derive from the public expression of the speech act  
2448 of promising, and the norms on the public expression of the speech act of  
2449 promising derive from (and are explanatorily downstream from) the norms

2450 on the mental state conveyed by that speech act.<sup>13</sup> If the norms on this mental  
 2451 state were different, the norms on the speech act would be different, too.<sup>14</sup>

2452 Because *every* promise expresses a resolution—that is, because every  
 2453 promise, regardless of whether the promisor is sincere or believed by the  
 2454 promisee, *conveys* that the promisor has resolved to act—*MSF* has broader  
 2455 applicability than do the theories discussed in the previous section, and can  
 2456 accommodate the marginal cases that these other theories cannot. When  
 2457 your highly unreliable friend tells you that he is going to pay you back the  
 2458 money you lend him, your expectations about repayment are not raised—but  
 2459 your friend does convey to you that he’s resolved to pay you back (even if you  
 2460 doubt that he’ll carry out this resolution). You can express resolutions in the  
 2461 absence of social conventions. A promise made to someone who later dies  
 2462 expresses a resolution in just the same way as a promise made to someone  
 2463 who remains living does. And one can convey that she has resolved to  $\Phi$   
 2464 without thereby communicating an intention to obligate herself. *MSF* is  
 2465 broader than the views discussed in the previous section because it derives  
 2466 promissory obligation from a core component of the speech act of promising  
 2467 itself, rather than from the content of particular promises or the various ways  
 2468 in which a promise interacts with the world (e.g., by creating expectations in  
 2469 the promisee, or being part of a social practice involving sanctioning).

- 
- 13 The order of explanation goes from mental state to speech act rather than vice-versa, for there is a clear sense in which mental states are independent of and prior to the speech acts that express them, and in which speech acts are not independent of and prior to the mental states they express. One can properly and without any insincerity be in a given mental state without expressing it via a speech act. Since the mental states can properly function entirely independent of the speech acts, it would be strange to derive the norms on them from the norms on speech acts. But speech acts are not independent from mental states in the same way: they express that the agent is in a particular mental state, and it is problematically insincere to have the speech act without that mental state. This dependency makes it quite natural to derive the norms on the speech act from the norms on the mental state.
- 14 *MSF* states that you can obligate yourself at will (pending the promisee’s acceptance) by publicly expressing that you have formed a conditional resolution to act. Does this make it a normative power view of the sort discussed in section 1? It doesn’t, because all theories of promissory obligation grant that promising involves deliberately creating a new obligation by doing something intentional (e.g., raising someone’s expectations, or participating in a social practice, or conveying a certain intention). What makes normative power views distinct is the explanatory structure that they take (i.e., a transcendental argument grounded in the legitimate interest we have in being able to alter the normative situation by fiat). *MSF* proposes a different explanatory structure: the basis of the promissory obligation is underlying the norms on the mental state conveyed by the promise, and what the public expression of that state commits one to. Thanks to an anonymous referee for discussion of this point.

2470 So far, I have illustrated what **MSF** can accomplish if it is true. My aim in  
 2471 the rest of this paper is to illustrate how plausible **MSF** in fact is. I begin by  
 2472 defending the claims that promises express resolutions (claim 1). I defend  
 2473 claims (2) and (3) in subsequent sections.

### 2473 **Defending Claim (1): Promises Express Resolutions**

2475 What do you express or convey when you perform the communicative act  
 2476 that constitutes making a promise? That is, what mental state do you at least  
 2477 purport to be in when promising, and are you really in if your promise is  
 2478 sincere?<sup>15</sup> A natural idea is that promising to  $\Phi$  expresses that you plan to  
 2479  $\Phi$ , and are serious about carrying out this plan. Someone who promises to  $\Phi$   
 2480 communicates that she really will  $\Phi$ , even if she doesn't feel like doing so at  
 2481 the time, or a better option arises, etc. When I promise you that I will attend  
 2482 your show tomorrow, I'm telling you that I am going to be there, even if faced  
 2483 with barriers to action that might otherwise prevent me from going. If my  
 2484 promise is sincere, I really do have such a plan. If I am insincere, I express  
 2485 to you that I am serious about going to your show without actually being so  
 2486 committed.

2487 How can we cash out this notion of a serious plan? We cannot appeal  
 2488 merely to desires, for desires aren't normatively committing in the way that  
 2489 promises are: promising to  $\Phi$  pro tanto morally obligates you to  $\Phi$ , but desiring  
 2490 to  $\Phi$  clearly does not. Moreover, we often sincerely promise to do things that  
 2491 we don't desire to do; you can sincerely promise your department chair that  
 2492 you'll attend the next faculty meeting, even though this is not how you desire  
 2493 to spend your Friday afternoon.<sup>16</sup> Intentions are a mental state with more  
 2494 stability than desires, as Michael Bratman (1987, 18–20) and others have  
 2495 argued. If I intend to go to your band's show tonight, then I have settled the  
 2496 matter of what I am going to do. I should not continue to deliberate about  
 2497 it or revise my plans for no good reason. Similarly, promising you that I will  
 2498 attend the show settles the matter of what I am going to do. It would be  
 2499 inappropriate to continue deliberating about what to do or to revise my plans

---

15 I refer not to sincerity in the sense of being well-meaning or earnest, but to *communicative sincerity*. When A utters a speech act S that expresses state M, A is communicatively sincere if and only if A really is in state M.

16 You might have other desires that would be satisfied by going to the meeting (e.g., getting on the chair's good side). But we can imagine scenarios in which this is not the case.

2500 unilaterally. Intentions might therefore be a decent candidate for the mental  
 2501 state expressed by promising.

2502 However, mere intentions are not stable enough to capture the seriousness  
 2503 of the plan and the strong sense of commitment involved in promising. Sup-  
 2504 pose that I intend to go to your band's show tonight because I have nothing  
 2505 better to do. You attempt to solicit a promise from me to attend the show. I  
 2506 say, "I promise to be there. And as of now, I intend to go. But you should  
 2507 know that I'm not committed to refraining from revising that intention. My  
 2508 plans might change between now and then, especially if I get a better offer."  
 2509 Such a statement does not seem to be a genuine promise. This is because it is  
 2510 unproblematic for me to intend to go to your band's show tonight and then  
 2511 abandon my intention because a more appealing offer comes along. But such  
 2512 circumstances would not license my breaking a promise to you to see your  
 2513 band play.

2514 This shows us that promises express something stronger and with more  
 2515 stability than typical intentions. Consider the special kind of intentions that  
 2516 we tend to form at the start of a new year—what we often call *resolutions*. These  
 2517 are particularly serious and stable intentions that we're strongly committed  
 2518 to, usually about important goals that we expect might be very difficult to  
 2519 attain. I take resolutions to be intentions that are especially robust or resistant  
 2520 to revision.<sup>17</sup> Resolutions are necessary when you plan to act and care about  
 2521 whether you do so but suspect that some temptation or other barrier to action  
 2522 (such as laziness, aversion to an unpleasant task, apathy, etc.) might cause you  
 2523 to abandon your plan were you not to bolster it somehow. Forming a resolution  
 2524 is one means by which you can bolster your plan and more effectively resist  
 2525 temptation.<sup>18</sup> It is plausible that when I promise to  $\Phi$ , I express that I have  
 2526 resolved to  $\Phi$ —that is, that I plan to  $\Phi$ , and that I care enough about whether  
 2527 I do so that I will not reconsider or abandon that plan, even in the face  
 2528 of temptation, laziness, better offers, and the like. Unlike intentions, it is

---

17 The details of how we cash out resolutions do not matter, so long as there is some coherent notion of resolution that implies that we are irrational if we over-hastily revise or fail to act on our resolutions without a good excuse. In Liberman (2016), I argue that resolutions consist in an intention to act coupled with the desire not to reconsider that intention, and offer an objection to Richard Holton's (2009) closely related view according to which a resolution is an intention to act coupled with an intention not to reconsider.

18 I don't mean to claim that resolutions are the *only* effective means of resisting temptation. There are other means by which you can resist temptation, which can be more effective than resolution-making; the best way to refrain from the temptation to drink tonight might not be to form a resolution, but to lock the liquor cabinet and give the key to a reliable friend.

2529 problematic to abandon a resolution because you no longer feel like acting,  
2530 or because a better offer comes along. So a resolution-based account will  
2531 not overgenerate cases of permissible promise-breaking in the way that an  
2532 intention-based account would.

2533 However, not every publicly expressed resolution counts as a promise; to  
2534 simply announce in your presence that I have resolved to run a marathon  
2535 is not to promise you that I will do so. This is because promises require a  
2536 second party; the acceptance or uptake of the promisee is essential to making  
2537 a promise. Promissory resolutions must therefore take account of promisee  
2538 acceptance. They can do this if they are *conditional on* the acceptance of the  
2539 promisee, in the way my resolution to go for a picnic tomorrow might be  
2540 conditional on the weather being good. If I offer to promise to run a marathon  
2541 with you, I convey that I resolve to run *on the condition that you accept* (and do  
2542 not subsequently reject) my promissory offer. A valid promise is successfully  
2543 created (and a pro tanto moral obligation generated) only if you accept my  
2544 offer.

2545 I propose that promises express a resolution to act, conditional on the  
2546 acceptance of the promisee. However, we might worry that this account over-  
2547 generates cases of legitimate promise-making. For we can imagine publicly  
2548 proclaimed resolutions that are explicitly conditional on someone else's accep-  
2549 tance or agreement, but that do not seem to count as promises. For example,  
2550 suppose I tell my personal trainer that I resolve to lift weights with her five  
2551 times a week, but only on the condition that she agree to work with me; she  
2552 agrees. I have announced that I have resolved to train five days a week, condi-  
2553 tional on her acceptance of this resolution, and she has accepted. But I don't  
2554 seem to have made her a promise. How do we distinguish genuine promises  
2555 from announcements that one has a conditional resolution?

2556 The response to this worry is simple: it is plausible that a necessary precon-  
2557 dition for validly promising is recognizing that you are making a promise in  
2558 the first place. In general, successfully engaging in an intentional action that  
2559 alters the normative situation requires a recognition of what you are doing,  
2560 e.g., you must recognize that you are granting consent in order to successfully  
2561 do so, and must be aware of the fact that you are transferring property in order  
2562 to make a gift, etc. Similarly, it is plausible that I need to understand that I  
2563 am making a promise in order to successfully do so. Likewise, it's plausible  
2564 that the promisee must be aware of the fact that she is accepting a promise, as  
2565 well. If the trainer does not *take herself to be accepting a promise* and thereby  
2566 generating a valid promissory obligation when she says she'll work with me,

2567 then what she is doing is agreeing with a resolution, and not accepting a  
2568 promise.<sup>19</sup>

## 2569 **4 Defending Claim (2): Resolutions Commit You to Acting**

### 4201 *Resolutions and Rational Commitment*

2571 Someone who resolves to  $\Phi$  incurs a self-imposed, pro tanto, subjective rational  
2572 obligation to  $\Phi$ . I will refer to such obligations as *rational commitments*.<sup>20</sup>  
2573 The easiest way to get a sense of what I mean by rational commitment is to  
2574 consider the way in which holding one belief can commit you to holding  
2575 another. Philosophers frequently talk about the ways in which our beliefs  
2576 commit us, e.g., because Jack believes that only consequences are relevant for  
2577 moral assessment, he is committed to believing that he ought to kill one person  
2578 to save two. Such commitments can come apart from what you objectively  
2579 ought to believe, all-things-considered: perhaps Jack should not believe that  
2580 he ought to kill one to save two. A resolution commits you to acting in a

---

19 MSF proposes that promises are best understood as conditional resolutions, but we must appeal to the concept of a promise to distinguish which resolutions are promise-generating and which aren't. Is this problematically circular? It isn't, for MSF is not meant to be a descriptive account of what promise-making consists in or a tool for identifying which utterances count as promises and which don't. Rather, MSF is an account of the nature of promissory obligation and of the normative force of promises. MSF posits that this ultimately stems from the norms on the conditional resolution that is publicly expressed when you make a promise. Promisor and promisee must both take themselves to be participating in promise-making in order to determine which resolutions will play this role. But this awareness of what they are doing is not the fundamental normative mechanism and does not provide any deep explanation of how and why promises bind. And in general, it is unproblematic to appeal to the concept of X as one part of an explanation of the nature or normative force of X. For example, suppose I am offering an account of the nature and normative force of giving, according to which A gives X to B (in a moral, rather than a legal sense) if and only if A intends to transfer X to B in an irrevocable way. Some irrevocable transfers will fail to count as gifts—say, those that are made under duress and are perceived by A and B as threats. We can unproblematically state that A must intentionally conceive of themselves as giving X to B in order for X to count as a gift, and that gift-giving generally involves conceiving of oneself as making a gift vs. acceding to a threat. This doesn't diminish the explanatory force of the account of giving as irrevocable transfer; appealing to the *concept* of a gift is part of what cashing out the underlying *nature* and *normative force* of giving requires. Something similar is true for my account of promissory obligation: we must appeal to the concept of a promise to cash out the underlying nature and normative force of promising. Thanks to an anonymous referee for discussion of this point.

20 For more on the concept of commitments (as distinct from reasons and all-things-considered obligations), see Shpall (2013, 2014), as well as Liberman and Schroeder (2016).

2581 similar way as believing that  $p$  commits you to believing the obvious and  
 2582 relevant consequences of  $p$ : someone who believes that  $p$  and that  $p$  entails  $q$   
 2583 but does not believe that  $q$  when the question of whether  $q$  is salient fails to  
 2584 act on her rational commitments, and as a result her overall set of beliefs is  
 2585 not as complete and coherent as it should be. Similarly, someone who resolves  
 2586 to  $\Phi$  at a particular time and then fails to intend to do so at that time because  
 2587 she readily abandons this resolution fails to act on her rational commitments,  
 2588 and as a result is not as effective a planning agent as she could be.<sup>21</sup> The  
 2589 rational obligation to act on one's resolutions stems from a broad demand for  
 2590 coherence in one's long-term plans; we are rationally committed to acting on  
 2591 our resolutions because this is essential for effectively carrying out our plans  
 2592 and acting in line with our important goals and values in the long term. To  
 2593 be clear, the obligation to fulfill your resolutions is not the same as a narrow  
 2594 requirement to be instrumentally rational or risk incoherence.<sup>22</sup> Rather, it is  
 2595 grounded in a broader demand for a more holistic sort of coherent planning  
 2596 agency. In order to meet our most important and difficult long-term goals, we  
 2597 must bolster ourselves against succumbing to temptation in ways that would  
 2598 undermine these goals. Resolution-keeping enables us to do this.

2599 We can best illustrate how resolutions incur rational commitments—that is,  
 2600 how they impose subjective, pro tanto, rational obligations—with an example.  
 2601 Suppose you are generally hesitant to try new foods, and are deeply entrenched  
 2602 in the habit of eating pizza for lunch every day. You very much want to expand  
 2603 your culinary horizons, but are such a creature of habit that you are unlikely  
 2604 to do so unless you force yourself into it somehow. So you resolve to go to a  
 2605 Thai restaurant for lunch today, knowing that if you don't form this resolution  
 2606 you are likely to fall back into your pizza habit. When you leave your office to  
 2607 go eat lunch, you abandon your plan to go to the Thai place and head to the  
 2608 pizzeria instead, deciding that you might as well just eat pizza today, since  
 2609 it's easier for you to order from a familiar menu. Something is wrong with  
 2610 you in this picture; resolving to eat Thai food in order to expand your culinary  
 2611 horizons and then changing your mind without good reason is problematically

21 Time indexing is necessary to avoid over-generating cases: if I resolve to go to eat Thai food at noon, it's no problem that I haven't yet formed the intention to do so at 9 AM. But it would be problematic if lunchtime rolls around and I haven't formed such an intention.

22 For influential treatments of conditional normative requirements grounded in coherence as the basis of instrumental rationality, see Broome (1999), who articulates a wide-scope view, and Kolodny (2005), who articulates a narrow-scope view; see also Way (2010) for an overview of this debate.

2612 self-undermining. It is irrational or incoherent to resolve to eat Thai food  
 2613 because you care about broadening your horizons, and then decide to stick  
 2614 with your pizza habit because it is easier. It's not that eating pizza every day is  
 2615 *independently* irrational; having pizza all the time is rationally permissible,  
 2616 if boring. Rather, it's that it is irrational to *resolve* to eat Thai food, and then  
 2617 abandon this resolution for no good reason.

2618 Compare abandoning a resolution to eat Thai food to merely *desiring* to  
 2619 have Thai food and eating pizza instead. This is perfectly acceptable; failing  
 2620 to act on a particular desire is not irrational. Or compare it to the case in  
 2621 which you *intend* to have Thai food but don't really care about whether this  
 2622 plan changes. You do not display any irrationality if you change your mind  
 2623 because you suddenly have a craving for a sandwich; you've simply made a  
 2624 permissible change of intention on the basis of a change in desire. Changing  
 2625 your resolution on a similar basis is not so innocuous. When you abandon  
 2626 your resolution to go to the Thai restaurant for lunch, you are undermining  
 2627 your own endorsed goals and plans: you value culinary diversity, and adopt  
 2628 this as one of your goals, but do not succeed in attaining it.

2629 We can best see how abandoning a resolution without good reason is irra-  
 2630 tional by comparing a pair of similar cases. Suppose that you and I each have  
 2631 a reason of strength X to eat Thai food and expand our respective culinary  
 2632 horizons. And suppose we each care about expanding our culinary horizons  
 2633 to the same extent. But only I take action about it: I *resolve* to eat Thai food,  
 2634 making it a part of my plan and adopting it as one of the concerns that I will  
 2635 focus on. If I change my mind and have pizza for lunch, I go wrong in a way  
 2636 that you do not when you have pizza. This is not to say that it is *never* permis-  
 2637 sible to abandon a resolution. Sometimes, there are weightier considerations  
 2638 in favor of revising a resolution than there are in favor of maintaining it—say,  
 2639 if you resolve to have Thai food and an old friend unexpectedly offers to meet  
 2640 you for lunch at the pizza place, or you realize that Thai food often contains  
 2641 ingredients to which you are allergic. However, resolutions rationally commit  
 2642 you to acting unless there is a good excuse in place.

#### 4<sup>42</sup> *Resolutions and Moral Obligation*

2644 It might seem *prima facie* puzzling to claim that the norms on promises can be  
 2645 explained by appeal to the norms on resolutions. I've argued in the previous  
 2646 section that resolving *rationally* commits you to acting. But promising *morally*  
 2647 obligates you to act, as well. Promissory obligations are also directed towards



2648 another person in a way that most resolutions are not: a promise does not  
 2649 create a generic moral obligation to act in a certain way, but is owed *to the*  
 2650 *promisee*, who is uniquely wronged should the promise be broken. We can  
 2651 assess the directedness of an obligation by asking whether anyone is uniquely  
 2652 positioned to forgive (or has special standing to resent) the violator, where  
 2653 resentment is understood as a negative evaluative attitude that is appropriate  
 2654 only in response to a personal violation. For example, suppose a passenger  
 2655 on an airplane has a medical emergency. The flight crew asks if any medical  
 2656 professionals are on board. There is a doctor on board (of medicine, not of  
 2657 philosophy!) who could easily assist the passenger, but who fails to do so  
 2658 because she prefers not to miss any scene of her in-flight movie. The other  
 2659 passengers on the plane might criticize or blame the doctor for her callous  
 2660 behavior. But the distressed passenger has unique standing to resent or forgive  
 2661 the doctor in a distinctive way. The doctor owes it *to the distressed passenger*  
 2662 to help him; this is not a generic obligation to be beneficent, but an obligation  
 2663 owed to a particular person. How can a resolution yield a commitment that is  
 2664 both *moral* and *directed* in this sort of way?

2665 I have three responses to this important question. First, I want to ease the  
 2666 explanatory burden on *MSF*. I am not trying to offer a complete or exhaustive  
 2667 account of all of the ways in which promises can morally obligate us, which  
 2668 means that *MSF* does not need to explain all of our intuitions about the  
 2669 robust moral force of promises. Rather, there are often *multiple* sources of  
 2670 moral obligation to keep any given promise, and these all contribute to the  
 2671 overall moral force of the promise.<sup>23</sup> These sources are the morally important  
 2672 considerations raised by the other theories of promissory obligation discussed  
 2673 in section 1: I ought to keep my promise to you to pick you up from the  
 2674 airport because it will harm you if I don't, *and* because failing to do so would  
 2675 problematically free-ride on a valuable social practice, *and* because I have  
 2676 publicly conveyed to you that I have resolved to do so. This last source of  
 2677 obligation stems from *MSF*, and is present in all cases of promising.

2678 Consider a pair of examples to illustrate how the moral force of a promise  
 2679 can vary (and will be stronger when other considerations are in play). First,  
 2680 suppose you promise your mother that you will not sell a treasured family  
 2681 heirloom that is in your possession. Assume that there aren't any strong  
 2682 reasons for you to sell it; you're doing fine financially, and the heirloom is

23 When there are multiple sources of obligation to keep a promise, the obligation is over-determined, and the moral force of it is stronger; see the case about reasons not to sell a family heirloom below.

2683 not worth very much. The moral reasons you have to keep this promise are  
 2684 not exhausted by the fact that you've conveyed to your mother a resolution to  
 2685 hang on to the heirloom; they also stem from harm avoidance, family loyalty,  
 2686 the importance of maintaining family traditions, and the like. And so the  
 2687 moral force of the promise to refrain from selling the heirloom will be quite  
 2688 strong.

2689 In cases in which no other reasons to keep a promise are present—that is,  
 2690 in which no expectation is formed, in which no trust can be violated, etc.—we  
 2691 see that promises generate comparatively weak moral obligations. Second,  
 2692 suppose that we are airplane seatmates stranded on the tarmac because of  
 2693 a flight delay. I ask you to promise me that you will write a negative review  
 2694 about the airline's customer service when you get off the flight, and you make  
 2695 the promise (and thereby resolve to do so, conditional on my acceptance).  
 2696 Assume that there are no independent sources of moral obligation to keep this  
 2697 promise—we are strangers, I will never know whether you write the review,  
 2698 doing so will in no way affect the airline's business, etc. You nevertheless  
 2699 seem to be pro tanto morally obligated to write the review. If you have an  
 2700 opportunity to easily write a negative review and you fail to do so without a  
 2701 good excuse, you will not be entirely as you ought to be. But neither will you  
 2702 be failing morally in a drastic way, for your obligation to keep the promise  
 2703 to me is not an especially important one. *MSF* purports only to explain this  
 2704 type of obligation: the relatively weak moral force that stems from the bare  
 2705 act of promise-making itself, including in marginal cases. This mitigates  
 2706 the explanatory demands on *MSF*, which needs to explain only this form of  
 2707 relatively weak moral obligation.

2708 Second, the conditionality of the resolution on the interest of the promisee  
 2709 is what enables it to yield a directed obligation. Successful promise-making  
 2710 requires the acceptance of the promise by the promisee, and I've argued  
 2711 that promissory resolutions are conditional on such acceptance. Accepting  
 2712 a promise implies that you are interested in the promise being made, and  
 2713 (usually) that you are interested in it being kept.<sup>24</sup> It's plausible that resolutions  
 2714 can generally yield directed obligations when the interests of another person

---

24 In typical cases, the promisee accepts the promise because she wants the promised action to occur. In deviant cases, the promisee might accept the promise for other reasons; Vera Peetz (1977) describes a case in which you accept your neighbor's promise to give you homemade jam not because you actually want her jam, but because you do not want to hurt her feelings by declining her offer. In this case, you are interested in the promise's being made—for this is necessary to spare your neighbor's feelings—even though you do not care whether it is kept.

2715 are directly involved in the core content of the resolution. Another person's  
2716 interest in how you are going to act is normatively relevant, and in some  
2717 circumstances—including promissory resolutions—this can ground a directed  
2718 obligation.

2719 We can illustrate how another person's interest in a resolution can yield  
2720 a directed obligation by considering a series of cases. Suppose you resolve  
2721 in secret to mow your elderly neighbor's lawn as a favor to him. In doing so,  
2722 you incur a rational commitment. If you fail to keep the resolution without  
2723 a good excuse for changing your mind, you will have acted irrationally, but  
2724 will not have done anything morally problematic. If you inform your neigh-  
2725 bor that you have resolved to mow his lawn, you will also incur a directed,  
2726 expectation-based obligation to your neighbor to either mow the lawn or alter  
2727 his expectations. And if you neither mow the lawn nor warn him in advance  
2728 that you do not plan to do so, your neighbor is entitled to resent or forgive you  
2729 for violating this expectation-based obligation. So too would be any third party  
2730 who has a stake in the matter and whose expectations were similarly raised  
2731 by your announcement of your resolution; if your neighbor's landscaper was  
2732 present when you announced your resolution, she is likewise entitled to resent  
2733 you for neither mowing the lawn nor warning her that you've changed your  
2734 mind.

2735 However, suppose you inform your neighbor that you have formed a reso-  
2736 lution that is explicitly conditional on his interests—you say, "I've resolved to  
2737 mow your lawn this weekend, so long as you want me to." In this case, it is  
2738 somewhat plausible that you incur a directed obligation to your neighbor to  
2739 mow the lawn. For your resolution conveys that you plan to mow his lawn,  
2740 and that this plan depends not on *your* interests or desires, but on *his*. The  
2741 conditionality of this resolution places the plan in his hands. If you don't mow  
2742 the lawn and instead simply warn your neighbor ahead of time that you have  
2743 changed your mind, you have perhaps done something morally problematic.  
2744 It would not seem terribly out of place for your neighbor to resent or forgive  
2745 you, in a way that would be odd if the resolution were not conditional on  
2746 his interests. But if you warn the landscaper ahead of time that you will not  
2747 be carrying out the conditional resolution as your neighbor wants you to, it  
2748 *would* seem out of place for the landscaper to resent or forgive you; you have  
2749 discharged your expectation-based obligation to her, and owe nothing more.<sup>25</sup>

---

25 There might be cases in which it is morally inappropriate for independent reasons to give the landscaper a warning instead of mowing the lawn yourself—say, if you know that the landscaper

2750 It is even more plausible that a directed obligation is formed if your resolu-  
 2751 tion is explicitly conditional not on your neighbor's *desire* that you mow his  
 2752 lawn, but on his active *acceptance* of your resolution to do so—that is, if you  
 2753 say, “I’ve resolved to mow your lawn, but will only do so if you remain actively  
 2754 on board with this plan.” For in that case, you have conveyed to your neighbor  
 2755 that you plan to mow his lawn, and that this plan depends not on your whims  
 2756 or desires but on his active endorsement of the plan. If you change your mind  
 2757 and don’t mow the lawn in spite of his continued uptake, it seems appropriate  
 2758 for him to resent or forgive you, even if you do warn him in time. Again,  
 2759 the landscaper does not seem similarly positioned to resent or forgive you,  
 2760 since the resolution was not conditional on her agreement. Your resolution  
 2761 yields not just a generic obligation owed to anyone who overheard you, but  
 2762 an obligation that is directed specifically towards your neighbor. In general, if  
 2763 you resolve to  $\Phi$  conditional on the agreement, acceptance, or uptake of A, it  
 2764 seems that you owe a directed obligation to A to  $\Phi$ . Promissory resolutions  
 2765 are always directed in this sort of way, because they are always conditional on  
 2766 the acceptance of the promisee.

2767 Third, conditional resolutions generate moral obligations because resolving  
 2768 conditionally on A’s acceptance and then failing to act when this condition is  
 2769 met fails to take proper consideration of A’s interests. To fail to adequately  
 2770 account for another person’s interests when you are engaged in a direct in-  
 2771 teraction with them is an interpersonal, moral sort of failing. It betrays the  
 2772 wrong kind of attitude to take towards another person, and can be construed  
 2773 as a form of disrespect or a problematic lack of moral concern. For example,  
 2774 suppose I proclaim to my family that I’ve resolved to give a particular heirloom  
 2775 to my cousin A, so long as she agrees to take it. She agrees, but I change my  
 2776 mind and give it to a different cousin, B. Assume A has no independent claim  
 2777 over the heirloom that B lacks; had I not formed a conditional resolution, it  
 2778 would be permissible for me to give the heirloom to A or B. Since I expressed  
 2779 a conditional resolution to give the heirloom to A, I seem to be *slighting* or  
 2780 *wronging* A when I give the heirloom to B instead. For I am failing to take her  
 2781 interests into account as I should and to give them their proper weight. My  
 2782 expressing a resolution conditional on A’s interests entails that she must be  
 2783 given special consideration; in the case where I’m simply deciding between A  
 2784 and B without expressing any resolution, I don’t wrong or slight A by choosing

---

has rearranged her entire schedule because she thought she didn’t need to mow your neighbor’s lawn. But generally, expectation-based obligations can be satisfied either by acting or altering the expectations.

2785 B, even if A wants the heirloom more than B does. Promissory resolutions  
2786 are those in which you are morally required to take special consideration  
2787 of the promisee's interests. To fail to do so is to fail to properly respect the  
2788 promisee, just as failing to give the heirloom to A when she has accepted  
2789 my conditional resolution fails to properly respect her. Granted, this may not  
2790 be a very serious or significant moral failing. But as we saw above, **MSF** is  
2791 burdened with establishing the existence only of a relatively weak sort of  
2792 moral obligation.

2793 In the last two sections, I have offered arguments for two of the core claims  
2794 that the argument for **MSF** relies on: (1) that promises express resolutions; (2)  
2795 that resolutions rationally and morally commit you to acting, all else being  
2796 equal. These claims are interesting in themselves, as they help us understand  
2797 the relationship between promises and resolutions, as well as the nature of the  
2798 commitments incurred by resolutions. With the addition of a third claim—  
2799 that is, the **determination claim**—this relationship between promises and  
2800 resolutions becomes significantly more interesting, for it can generate a new  
2801 kind of theory of promissory obligation.

2802 At present, I cannot fully defend the **determination claim** or offer a deep  
2803 explanation of why it holds. However, I can offer evidence to suggest that the  
2804 **determination claim** is both plausible and theoretically fruitful, insofar as it  
2805 can help explain the norms on speech acts other than promising.

## 2806 **5 Defending Claim (3): Evidence for the Determination** 2807 **Claim**

2808 The third claim on which **MSF** depends is the **determination claim**, which  
2809 states that the norms governing the mental state expressed by a speech act  
2810 (at least partially) determine what the norms on that speech act are, in both  
2811 sincere and insincere cases. In other words, the determination claim tells us  
2812 that saying that you're in state X commits you to behaving as if you are in that  
2813 state, regardless of whether you in fact are.

2814 My first piece of evidence for the **determination claim** is the observation  
2815 that, in general, we need to assume the truth of something like the **deter-**  
2816 **mination claim** in order to have fair and productive social interactions. The  
2817 **determination claim** states that you must act as if you really are in a particular  
2818 mental state when you convey to others that you are in that state. People are  
2819 not mind-readers, and our conversations and social interactions are generally

2820 presumed to be cooperative. In light of this, it would be unfair to expect people  
 2821 to be able to ascertain when our utterances are sincere and when they are not.  
 2822 It follows that we must be able to take what people say at face value if we are  
 2823 to have productive interactions with them, at least in typical circumstances.<sup>26</sup>  
 2824 Doing so enables us to respond appropriately to them: to predict what they  
 2825 might do and say next, and to alter our behavior in light of theirs, etc.

2826 If we weren't entitled to presume that people really were in the states they  
 2827 purport to be in, we wouldn't be able to interact with them very effectively. If  
 2828 I promise to  $\Phi$ , my promisee must be entitled to presume that I am committed  
 2829 to  $\Phi$ ing, lest she be at a loss for how to respond to me. And because she cannot  
 2830 know my inner mental state, it would be unreasonable and unfair to expect  
 2831 anything else of her. Accordingly, my promisee is entitled to interact with me  
 2832 as if I have in fact resolved to  $\Phi$ —which might involve her believing that I plan  
 2833 to  $\Phi$ , structuring her future plans and behavior around my  $\Phi$ ing, or simply  
 2834 responding appropriately in the moment to my commitment to  $\Phi$ ing. And  
 2835 as promisor, I should in turn behave as if I really have resolved, and thereby  
 2836 enable my interlocutor to take my promissory utterance at face value—which  
 2837 is to say, I should act in accordance with the **determination claim**.

2838 The rest of my evidence for the **determination claim** is circumstantial: there  
 2839 exist other cases in which publicly conveying that you are in a mental state  
 2840 M by performing a speech act S commits you to acting in whatever way is  
 2841 required by M, regardless of whether you really are in M. Promising is not  
 2842 unique in this regard, but is part of a general pattern. This gives us good reason  
 2843 to think that the **determination claim** is broadly true, and is not merely an *ad*  
 2844 *hoc* principle that applies only to the case of promises and that I am invoking  
 2845 out of the blue to defend **MSF**.

---

26 Because the mental state is explanatorily prior to the speech act (see footnote 13), the mental state is also explanatorily prior to the action that results from the speech act. Suppose that I am in state X, and I convey this by saying "I'm in X." How others can reasonably expect me to act on the basis of this utterance is fundamentally determined not by my utterance that I'm in X, but by what *being in state X* commits me to. To be fair to others, I must act as if I am in state X, unless there is some reason not to take my utterance at face value—say, if it is mutually understood that I am in a strategic (e.g., game-playing) context.

5.4.1 *Forgiving*

2847 It is plausible that forgiving someone expresses that you have repudiated or  
2848 foresworn blaming them.<sup>27</sup> Repudiating blame of A for doing X plausibly  
2849 commits one to ceasing blaming A for X in the moment, and to refraining  
2850 from actively expressing blame towards A for X again in the future. Publicly  
2851 expressing forgiveness—regardless of whether you have privately foresworn  
2852 blame—likewise seems to commit you to refraining from expressing blame in  
2853 these same ways. That is, publicly conveying that you have repudiated blame  
2854 by performing a speech act of forgiveness commits you to behaving as if you  
2855 really have repudiated blame.

2856 To illustrate, suppose that Anna forgets that today is her wedding anniversary,  
2857 and she fails to meet her spouse Betty for a celebratory dinner they have  
2858 planned. Betty knows that she is likely to hold this mistake against Anna,  
2859 and doesn't want to damage their relationship by doing so. So she decides  
2860 to forswear blaming Anna for her oversight. This private forswearing of  
2861 blame commits her to refraining from expressing blame towards Anna. And  
2862 telling Anna that she has forgiven her likewise commits Betty to refraining  
2863 from expressing blame towards Anna. This is so even if Betty is insincere, and  
2864 forgives Anna not because she has in fact foresworn blaming her, but because  
2865 she wants to avoid conflict. Betty's interpersonal utterance of forgiveness  
2866 nevertheless commits her to refraining from openly expressing blame. As  
2867 with promising, the public expression of a forswearing of blame changes the  
2868 nature of the norms to which you are subject; it transforms a private commit-  
2869 ment to refraining from engaging in blaming activities into an interpersonal,  
2870 directed obligation to avoid continuing to blame the wrongdoer in the future.  
2871 Finally, and again as with promising, this is compatible with the existence of  
2872 additional explanations of why you should refrain from expressing blame in  
2873 particular cases; perhaps Betty has an obligation to refrain from expressing  
2874 blame because she has publicly expressed that she has forgiven Alice, *and*  
2875 because she is independently obligated to avoid acting unfairly, and expressing  
2876 blame would be unfair since Betty forgot their anniversary last year.

---

27 This is a common view of forgiveness in both everyday practice and the philosophical literature; the view is usually attributed first to Bishop Butler. See discussion of this view in Griswold (2007).

## 5.2 *Apologizing*

2878 Similarly, apologizing for  $\Phi$ ing expresses that you regret or are sorry that you  
 2879  $\Phi$ ed. When someone regrets an action in this way, they incur a commitment  
 2880 to actively taking responsibility for it somehow. What this involves will vary  
 2881 in different cases; a cheating spouse who regrets their infidelity is obligated  
 2882 to avoid straying again, while a party guest who regrets spilling red wine on  
 2883 a white rug incurs a commitment to clean up the spill. Philosophical and  
 2884 popular consensus is that someone who publicly apologizes for  $\Phi$ ing likewise  
 2885 commits themselves to taking responsibility for their action.<sup>28</sup> And this is so  
 2886 regardless of whether the apology was sincere; publicly conveying that you  
 2887 regret your action by apologizing commits you to behaving as if you really  
 2888 have regretted your action.

2889 For example, suppose Christa catches her student Danny using his phone  
 2890 during class, in violation of her policy. Danny apologizes to Christa, which  
 2891 commits him to taking responsibility for his error (e.g., by admitting that he  
 2892 was wrong and refraining from using his phone in class again). This is so even  
 2893 if Danny is insincere, and apologizes only because he fears that Christa will  
 2894 dock his participation grade if he does not. As with promising and forgive-  
 2895 ness, a public apology transforms a personal feeling of regret that privately  
 2896 commits you to making amends into an interpersonal demand to make such  
 2897 amends. There may also be additional moral considerations present that re-  
 2898 quire agents to take active responsibility for their wrongdoing. But even if  
 2899 these considerations are not present, the mere (sincere or insincere) expres-  
 2900 sion of regret commits you—at least in a weak way—to taking responsibility  
 2901 for your wrongdoing.

## 5.3 *Asserting*

2903 There is much disagreement about what mental state assertion expresses and  
 2904 whether there are norms of further commitment on assertions; adjudicating  
 2905 between these views is too large a project to be adequately handled here. But  
 2906 we can assume a particular view to illustrate how assertion might pattern

---

28 For example, Mihaela Mihai (2013) notes that while philosophical accounts of apology vary, “there is a growing consensus that an authentic apology implies an acknowledgement that the incident in question did in fact occur and that it was inappropriate, a recognition of responsibility for the act, the expression of an attitude of regret and a feeling of remorse, and the declaration of an intention to refrain from similar acts in the future.”



2907 with promising, forgiving, and apologizing. Suppose for the sake of argument  
2908 that assertion expresses belief (which is compatible with assertion expressing  
2909 something else that includes belief as a component, such as justified belief  
2910 or knowledge). If you believe that  $p$  and are questioned about whether  $p$ , it  
2911 is plausible that you are normatively committed to defending or justifying  
2912  $p$ . It is also plausible that this same commitment is inherited by assertions:  
2913 some philosophers argue that someone who asserts that  $p$  takes on a special  
2914 commitment to the truth of  $p$ , which can be cashed out as a commitment  
2915 to justify or defend  $p$  to one's interlocutors.<sup>29</sup> This is plausibly so even in  
2916 insincere cases; if you are not going to retract your insincere assertion that  $p$ ,  
2917 you should be prepared to defend or justify it. If this is the case, then assertion  
2918 patterns with the other speech act/mental state pairs we've been discussing:  
2919 publicly conveying that you believe  $p$  by asserting that  $p$  commits you to  
2920 behaving as if you really believe  $p$ , and transforms your private commitment  
2921 to behave as if  $p$  is true into a commitment to defending  $p$  publicly.

## 5.2.4 Conclusion

2923 The **determination claim** is a natural and plausible explanation of the pattern  
2924 outlined above, for two reasons. First, speech acts are dependent on the mental  
2925 states they express, in a way that mental states are not dependent on the speech  
2926 acts used to express them. You can properly and without any insincerity be  
2927 in a particular mental state without expressing it via a speech act; it is no  
2928 problem to resolve to act without promising that you will do so. But all of the  
2929 speech acts we have been discussing express that the agent is in a particular  
2930 mental state, and it is communicatively insincere to perform the speech act  
2931 without being in that mental state. The direction of explanation proposed by  
2932 the **determination claim** tracks this dependency; it would not make sense for  
2933 the direction of explanation to go in the other direction.

2934 Because we have reason to think the **determination claim** is true—and  
2935 we also have reason to think that promises express resolutions, and that  
2936 resolutions rationally (and sometimes morally) commit us to acting—we have  
2937 reason to think that **MSF** is a viable theory of promissory obligation. And  
2938 **MSF** has appealing explanatory advantages. Since the mental state expressed  
2939 by a promise is both distinctive of promises and present in all cases, **MSF**

---

29 See, among others, Peirce (1935), Searle (1969), Brandom (1983), Wright (1992), Watson (2004), and MacFarlane (2005).

2940 captures the minimal, essentially promissory obligation that is always there,  
 2941 while being open to the pluralist idea that other theories can explain why we  
 2942 have stronger moral reasons to keep our promises in many cases. Moreover,  
 2943 the success of MSF should lead us to be optimistic about the possibility of  
 2944 providing similarly structured and equally resourceful accounts cashing out  
 2945 the norms on other sorts of speech acts in terms of their underlying mental  
 2946 states.\*

2947 Alida Liberman

2948  0000-0002-5182-569X

2949 Southern Methodist University

2950 aliberman@smu.edu

## 2951 References

- 2952 BRANDOM, Robert B. 1983. "Asserting." *Noûs* 17(4): 637–650, doi:10.2307/2215086.
- 2953 BRATMAN, Michael E. 1987. *Intentions, Plans and Practical Reason*. Cambridge, Mas-  
 2954 sachusetts: Harvard University Press.
- 2955 BROOME, John A. 1999. "Normative Requirements." *Ratio* 12(4): 398–419, doi:10.111  
 2956 1/1467-9329.00101.
- 2957 CHOLBI, Michael J. 2002. "A Contractualist Account of Promising." *The Southern*  
 2958 *Journal of Philosophy* 40(4): 475–491, doi:10.1111/j.2041-6962.2002.tb01913.x.
- 2959 COTTINGHAM, John G. 1985. "Review of Robins (1984)." *The Philosophical Quarterly*  
 2960 35(140): 315–318, doi:10.2307/2218913.
- 2961 DOWNIE, R. S. 1985. "Three Accounts of Promising." *The Philosophical Quarterly*  
 2962 35(140): 259–271, doi:10.2307/2218905.
- 2963 GRISWOLD, Charles L. 2007. *Forgiveness: A Philosophical Exploration*. Cambridge:  
 2964 Cambridge University Press, doi:10.1017/cbo9780511619168.
- 2965 HEUER, Ulrike. 2012. "Promising – Part 1." *Philosophy Compass* 7(12): 832–841, doi:10  
 2966 .1111/j.1747-9991.2012.00524.x.
- 2967 HOLTON, Richard. 2009. *Willing, Wanting, Waiting*. Oxford: Oxford University Press,  
 2968 doi:10.1093/acprof:oso/9780199214570.001.0001.

---

\* This paper is adapted from my dissertation, *The Mental States Theory of Promising* (2015), com-  
 pleted at the University of Southern California. More people than I can list here have helped  
 me to articulate and refine these ideas over the years. Special thanks are due to my committee  
 members Steve Finlay, Jake Ross, and Gary Watson, to my graduate school classmates (especially  
 Josh Crabill, N.G. Laskowski, Michael Milona, Caleb Perl, and Abelard Podgorski), to audiences  
 at USC and Bowling Green State University, and most importantly to my dissertation supervisor  
 Mark Schroeder.

- 2969 HOOKER, Brad. 2011. "Promises and Rule-Consequentialism." in *Promises and Agree-*  
 2970 *ments: Philosophical Essays*, edited by Hanoach SHEINMAN, pp. 235–252. Oxford:  
 2971 Oxford University Press, doi:10.1093/acprof:oso/9780195377958.003.0010.
- 2972 HUME, David. 1739. *A Treatise of Human Nature*. London: John Noon, at the White-  
 2973 Hart.
- 2974 —. 1978. *A Treatise of Human Nature*. 2nd ed. Oxford: Oxford University Press. Edited  
 2975 by Lewis Amherst Selby-Bigge and Peter Harold Nidditch, doi:10.1093/oseo/insta  
 2976 nce.00046221.
- 2977 —. 2007. *A Treatise of Human Nature. A Critical Edition. Volume 1: Texts*. 2nd ed.  
 2978 Oxford: Oxford University Press. Edited by David Fate and Mary J. Norton.
- 2979 KOLODNY, Niko. 2005. "Why Be Rational?" *Mind* 114(455): 509–563, doi:10.1093/mi  
 2980 nd/fzi509.
- 2981 KOLODNY, Niko and WALLACE, Richard Jay. 2003. "Promises and Practices Revisited."  
 2982 *Philosophy & Public Affairs* 31(2): 119–154, doi:10.1111/j.1088-4963.2003.00119.x.
- 2983 LEMOS, Noah M. 1987. "Review of Robins (1984)." *Philosophy and Phenomenological*  
 2984 *Research* 47(4): 685–688, doi:10.2307/2107242.
- 2985 LIBERMAN, Alida. 2015. "The Mental States First Theory of Promising." PhD disserta-  
 2986 tion, Los Angeles, California: Philosophy Department, University of Southern Cal-  
 2987 ifornia, <http://digitallibrary.usc.edu/cdm/ref/collection/p15799coll3/id/591683>.
- 2988 —. 2016. "Reconsidering Resolutions." *Journal of Ethics and Social Philosophy* 10(2),  
 2989 doi:10.26556/jesp.v10i2.98.
- 2990 —. 2020. "Consequentialism and Promises." in *The Oxford Handbook of Consequential-*  
 2991 *ism*, edited by Douglas W. PORTMORE, pp. 289–309. Oxford Handbooks. Oxford:  
 2992 Oxford University Press, doi:10.1093/oxfordhb/9780190905323.013.16.
- 2993 LIBERMAN, Alida and SCHROEDER, Mark. 2016. "Commitment: Worth the Weight." in  
 2994 *Weighing Reasons*, edited by Errol LORD and Barry MAGUIRE, pp. 104–120. Oxford:  
 2995 Oxford University Press, doi:10.1093/acprof:oso/9780199315192.003.0005.
- 2996 MACFARLANE, John. 2005. "Making Sense of Relative Truth." *Proceedings of the Aris-*  
 2997 *totelian Society* 105: 305–323, doi:10.1111/j.0066-7373.2004.00116.x.
- 2998 MARUŠIĆ, Berislav. 2013. "Promising Against the Evidence." *Ethics* 123(2): 292–317,  
 2999 doi:10.1086/668704.
- 3000 MIHAI, Mihaela. 2013. "Apology." in *Internet Encyclopedia of Philosophy*. University of  
 3001 Tennessee at Martin, <https://iep.utm.edu/apology/>.
- 3002 OWENS, David. 2012. *Shaping the Normative Landscape*. Oxford: Oxford University  
 3003 Press, doi:10.1093/acprof:oso/9780199691500.001.0001.
- 3004 PEETZ, Vera. 1977. "Promises and Threats." *Mind* 86(344): 578–581, doi:10.1093/mind  
 3005 /LXXXVI.344.578.
- 3006 PEIRCE, Charles Sanders. 1935. *Collected Papers of Charles Sanders Peirce, vol. 5: Prag-*  
 3007 *matism and Pragmaticism*. Cambridge, Massachusetts: Harvard University Press.  
 3008 Edited by Charles Hartshorne and Paul Weiss.

- 3009 PINK, Thomas. 2009. "Promising and Obligation." in *Philosophical Perspectives 23:*  
3010 *Ethics*, edited by John HAWTHORNE, pp. 389–420. Hoboken, New Jersey: John  
3011 Wiley and Sons, Inc., doi:10.1111/j.1520-8583.2009.00177.x.
- 3012 RAWLS, John. 1971. *A Theory of Justice*. Cambridge, Massachusetts: Harvard University  
3013 Press. Revised edition: Rawls (1999), doi:10.2307/j.ctvjf9z6v.
- 3014 —. 1999. *A Theory of Justice*. Cambridge, Massachusetts: Harvard University Press,  
3015 doi:10.2307/j.ctvkjb25m.
- 3016 RAZ, Joseph. 1977. "Promises and Obligations." in *Law, Morality and Society. Essays in*  
3017 *Honour of H.L.A. Hart*, edited by Peter M. S. HACKER and Joseph RAZ, pp. 210–228.  
3018 Oxford: Oxford University Press.
- 3019 ROBINS, Michael H. 1984. *Promising, Intending and Moral Autonomy*. Cambridge:  
3020 Cambridge University Press.
- 3021 SCANLON, Thomas Michael. 1990. "Promises and Practices." *Philosophy & Public*  
3022 *Affairs* 19(3): 199–226.
- 3023 —. 1998. *What We Owe to Each Other*. Cambridge, Massachusetts: Harvard University  
3024 Press.
- 3025 SEARLE, John R. 1969. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge:  
3026 Cambridge University Press, doi:10.1017/cbo9781139173438.
- 3027 SHEINMAN, Hanoch. 2008. "Promise as Practice Reason." *Acta Analytica* 23(4):  
3028 287–318, doi:10.1007/s12136-008-0033-1.
- 3029 SHIFFRIN, Seana Valentine. 2008. "Promising, Intimate Relationships, and Convention-  
3030 alism." *The Philosophical Review* 117(4): 481–524, doi:10.1215/00318108-2008-014.
- 3031 SHPALL, Samuel. 2013. "Wide and Narrow Scope." *Philosophical Studies* 163(3):  
3032 717–736, doi:10.1007/s11098-011-9841-z.
- 3033 —. 2014. "Moral and Rational Commitment." *Philosophy and Phenomenological Re-*  
3034 *search* 88(1): 146–172, doi:10.1111/j.1933-1592.2012.00618.x.
- 3035 SMITH, Holly M. 1987. "Review of Robins (1984)." *Noûs* 21(4): 604–608, doi:10.2307/  
3036 2215676.
- 3037 SOUTHWOOD, Nicholas and FRIEDRICH, Daniel. 2009. "Promises Beyond Assurance." *Philosophical Studies* 144(2): 261–280, doi:10.1007/s11098-008-9209-1.
- 3038 THOMSON, Judith Jarvis. 1990. *The Realm of Rights*. Cambridge, Massachusetts: Har-  
3039 vard University Press.
- 3040 WATSON, Gary. 2004. "Asserting and Promising." *Philosophical Studies* 117(1/2): 57–77,  
3042 doi:10.1023/B:PHIL.0000014525.93335.9e.
- 3043 WAY, Jonathan. 2010. "The Normativity of Rationality." *Philosophy Compass* 5(12):  
3044 1057–1068, doi:10.1111/j.1747-9991.2010.00357.x.
- 3045 WRIGHT, Crispin. 1992. *Truth and Objectivity*. Cambridge, Massachusetts: Harvard  
3046 University Press, doi:10.2307/j.ctv1rr6cwg.

# A Puzzle About Parsimony

PETER FINOCCHIARO

In this paper, I argue for the instability of an increasingly popular position about how metaphysicians ought to regard parsimony. This instability is rooted in an unrecognized tension between two claims. First, we as metaphysicians ought to minimize the number of ontological kinds we posit. Second, it is not the case that we ought to minimize the number of ideological expressions we employ, especially when those expressions are of the same ideological kind (e.g., the compositional predicates ‘is a part of’ and ‘overlaps’). I argue that the two claims are in tension with one other. At the very least, minimizing the number of ontological kinds posited entails minimizing the number of expressions employed—more specifically, the “ontologically committing” predicates. But, plausibly, the tension runs deeper than that. I suggest that minimizing the number of ontological kinds just is a specific way of minimizing the number of ideological expressions employed in stating a theory. The two activities target the same aspect of reality, the world’s metaphysical structure. I end by evaluating three different responses to this puzzle. Ultimately, I suggest that metaphysicians should treat the minimization of the number of ideological expressions as more important than it currently is treated.

Parsimony is among the most prominent methodological considerations in metaphysics. Yet beneath the surface there lurks a puzzle. I will bring this puzzle about parsimony to light. As I will show, the puzzle highlights a conceptual tension between several prominent positions in metaphysics. I will then offer three responses to the puzzle. Each response faces unique challenges.

First, I will make some starting assumptions. These assumptions are not unassailable. But each is independently plausible and each has broad support amongst metaphysicians.

Parsimony has traditionally been restricted to ontology: do not multiply entities beyond necessity. Lately, however, metaphysicians have turned their attention toward *ideological parsimony*. Ideological parsimony, as I understand it, concerns the primitive (i.e., undefined) terminology used to state a theory. Recently, many philosophers ([Brenner 2017](#); [Cowling 2013](#); [Schaffer 2015](#);

3079 Sider 2011; Turner 2015) have defended the claim that both ontological and  
 3080 ideological parsimony make a theory more worthy of our endorsement. I will  
 3081 assume that they are right.

3082 I will also assume what is sometimes called a *realist* or *externalist* interpret-  
 3083 ation of ideology. Like an analogous interpretation of ontology, ideological  
 3084 externalism states that the quality of a theory's ideology is ultimately judged  
 3085 by the extent to which it corresponds to objective reality, i.e., the world's  
 3086 metaphysical structure.<sup>1</sup> (Ideological externalism can be contrasted with ide-  
 3087 ological internalism, which states that the quality of a theory's ideology is  
 3088 judged by details internal to the theoretic process—e.g., the intelligibility of  
 3089 the terminology employed.)

3090 I will also adopt the orthodox approach to meta-ontology, *neo-Quineanism*.  
 3091 According to neo-Quineanism, a theory's ontological commitments are deter-  
 3092 mined by what the theory quantifies over when regimented with a suitably  
 3093 perspicuous language.<sup>2</sup> Finally, I will focus on theories about the fundamental  
 3094 nature of the world. While there might be versions of this puzzle that extend  
 3095 to non-fundamental theories, I do not have much to say about them. That is in  
 3096 large part because I do not have much to say in general about the relationship  
 3097 between fundamental and non-fundamental theories.

3098 These assumptions help generate a puzzle, one that highlights a conceptual  
 3099 tension in how some metaphysicians understand the role of parsimony in  
 3100 theory choice. This tension has, until now, gone unrecognized. To bring out  
 3101 the tension, I identify in section 1 four approaches to parsimony that differ  
 3102 along two axes: ontology/ideology and quantitative/qualitative. We seem to  
 3103 have an intuitive grasp on these approaches and understand the differences  
 3104 between them, in particular the differences between ontological and ideologi-  
 3105 cal parsimony. But in section 2, I argue that qualitative ontological parsimony  
 3106 entails a restricted version of quantitative ideological parsimony. This is a  
 3107 surprising and worrisome puzzle. It is surprising because it goes against our  
 3108 intuitive grasp of parsimony. It is worrisome because it seems inconsistent  
 3109 with a popular position amongst metaphysicians—i.e., that greater qualitative  
 3110 ontological parsimony makes a theory more worthy of endorsement but it  
 3111 is not the case that greater quantitative ideological parsimony makes a the-  
 3112 ory more worthy of endorsement. I then suggest that the entailment is no  
 3113 coincidence; qualitative ontological parsimony may be conceptually distinct

1 I discuss ideological externalism (as well as ideological internalism) in greater depth in Finocchiaro (2021, 963–969). See, also, Cowling (2013, 3983) and Sider (2011, 13).

2 See, inter alia, Quine (1948); van Inwagen (1998); Lewis and Lewis (1970).

3114 from quantitative ideological parsimony, but the most sensible applications  
 3115 of them target the same feature of reality, the world's metaphysical structure.

3116 In section 3, I discuss three available responses to this puzzle. First, we  
 3117 could resist the puzzle by rejecting neo-Quineanism. Second, we could down-  
 3118 play the significance of the puzzle by offering a more nuanced application of  
 3119 parsimony. Finally, we could reevaluate the value of quantitative ideological  
 3120 parsimony as a theoretical virtue. Ultimately, I favor the third response. Meta-  
 3121 physicists should value quantitative ideological parsimony more than they  
 3122 currently do.

### 3121 **1 Four Different Approaches to Parsimony**

3124 Many metaphysicians think that parsimony should play a role in theory choice.  
 3125 They have cited parsimony in support of theories as wide-ranging as compo-  
 3126 sitional nihilism (Horgan and Potrč 2008), bundle theory (Paul 2017), materi-  
 3127 alism (Churchland 1984), and nominalism (Melia 2008).

3128 But such metaphysicians often differ in how they use parsimony. Even when  
 3129 restricted to the ontology of a fundamental theory, there are two importantly  
 3130 different approaches they take. Some (e.g., Nolan 1997) tend to prefer the  
 3131 theory that minimizes the number of entities posited. Others (e.g., Lewis  
 3132 1973) tend to prefer the theory that minimizes the number of *kinds* of entities.  
 3133 Following a convenient shorthand from Cowling (2013), I will name these  
 3134 two different approaches (NO-Parsimony) and (KO-Parsimony), respectively.

3135 I won't take a stand on which approach is best.<sup>3</sup> I simply note that even  
 3136 those inclined toward (NO-Parsimony) also tend to be inclined toward (KO-  
 3137 Parsimony). More generally, among the metaphysicians who care about parsim-  
 3138 ony at all, most of them accept (KO-Parsimony).<sup>4</sup>

3139 We can also consider the parsimony of a fundamental theory's ideology.  
 3140 David Lewis, for example, claims that modal realism enables us "to reduce the  
 3141 diversity of notions we must accept as primitive" (1986, 4). Theodore Sider  
 3142 argues that compositional nihilism "allows us to eliminate the extra-logical  
 3143 (or perhaps quasi-logical) notion of 'part' from our ideology" (2013, 239). Both  
 3144 modal realism and compositional nihilism are ideologically parsimonious.

---

3 For two defenses of different uses, see Lewis (1973) and Tallant (2013).

4 For instance, Nolan (1997, 330) says "I claim that not only ought we not multiply types of entities beyond necessity, but that we should also be concerned not to multiply the entities of each type more than is necessary."



3145 For Lewis and Sider, the ideological parsimony of their theories provides a  
3146 reason to endorse them.

3147 Just as with ontology, there are two importantly different approaches to  
3148 ideological parsimony. Metaphysicians may prefer the theory that minimizes  
3149 the total number of terms that are employed but undefined within the theory  
3150 (“bits of ideology”). Or they may prefer the theory that minimizes the number  
3151 of kinds of terms so employed (“ideological kinds”). Adopting another short-  
3152 hand from Cowling (2013), I will name these approaches (NI-Parsimony) and  
3153 (KI-Parsimony), respectively.<sup>5</sup>

3154 I should note that it’s not obvious how to individuate ideological kinds. (The  
3155 same could be said about ontological kinds.) Metaphysicians often rely on the  
3156 imprecise but intuitive method of individuation by topic. For instance, there  
3157 is an ideological kind corresponding to color. All color predicates like ‘blue,’  
3158 ‘periwinkle,’ and ‘Pantone 19-4052’ are of this kind, as are relational predicates  
3159 like ‘is more saturated than.’ There is also an ideological kind corresponding  
3160 to modality. Primitive modal operators, predicates like ‘possibly true’ and  
3161 ‘consistent,’ as well as primitive dispositional predicates like ‘fragile’ are of this  
3162 kind. There is much more worth saying about the individuation of ideological  
3163 kinds.<sup>6</sup> Yet I do not think that my main argument is affected by this issue. In  
3164 what follows I will stick to the intuitive understanding just sketched.

3165 Some metaphysicians may deny that our use of ideological parsimony can  
3166 be neatly divided into (NI-Parsimony) and (KI-Parsimony). Yet the distinction  
3167 seems intuitive enough and many think there is something to it (e.g., Cameron  
3168 2012, 18; Cowling 2013, 3897). In addition, there are intuitive reasons to favor  
3169 (KI-Parsimony) and reject (NI-Parsimony). For one, (NI-Parsimony) seems  
3170 to force us to make objectionably arbitrary decisions. (NI-Parsimony) recom-  
3171 mends that, all else being equal, we minimize the number of compositional  
3172 predicates in our ideology. What this recommendation precisely amounts to  
3173 will depend on the resolution of issues that are too large to address here.<sup>7</sup>  
3174 To see the worry, though, suppose that there are no other relevant considera-  
3175 tions regarding our choice of compositional ideology. (NI-Parsimony) then  
3176 recommends that we employ a minimal expressively adequate set of predi-  
3177 cates. For composition, this can be achieved by choosing one from among ‘is

5 Some characterize ideology as concerning the *concepts* employed in stating a theory. I prefer my *linguistic* characterization, for reasons I state in Finocchiaro (2021, 961–963).

6 I do say much more in Finocchiaro (2019a). See, also, Cowling (2013) and Lewis (1986).

7 For example, it depends in part on whether composition is classically extensional true (see Parsons 2014, 4).



3178 a part of,' 'is a proper part of,' and 'overlaps' (supplemented with identity).  
3179 We are then faced with an unsettling question: which of these three should  
3180 we choose? Each option is unsavory because they all seem to commit us to  
3181 an unreasonable view about the fundamental compositional structure of the  
3182 world. Each option also seems impossible to motivate—what could justify  
3183 choosing one over the other? These worries about arbitrariness disappear if we  
3184 reject (NI-Parsimony) in favor of (KI-Parsimony). Compositional predicates  
3185 are (plausibly) of the same ideological kind. So there is no methodological  
3186 pressure to arbitrarily choose one predicate over the others.<sup>8</sup>

3187 Many metaphysicians nowadays think that both ontological parsimony and  
3188 ideological parsimony should play a role in theory choice. Why? Historically,  
3189 parsimony-based considerations have been defended on non-alethic grounds:  
3190 an ideologically parsimonious theory might be easier to comprehend, or an  
3191 ontologically parsimonious theory might be more aesthetically pleasing. But  
3192 such defenses are less popular nowadays since they are seen as relying on  
3193 reasons that should be irrelevant to theory choice in metaphysics. Nowadays,  
3194 most metaphysicians who think that parsimony should play a role in theory  
3195 choice think so because they think parsimony is truth-conducive.<sup>9</sup> This con-  
3196 nection between parsimony and truth holds for both ontological parsimony  
3197 and ideological parsimony. According to ideological externalism, a more ideo-  
3198 logically parsimonious theory conveys a more simple—and therefore more  
3199 likely to be true—picture of the world's structure. Yet metaphysicians are  
3200 less willing to extend this defense to quantitative ideological parsimony. (In-  
3201 tuitively, a theory that employs only 'is a part of' is not any more likely to  
3202 be true than a theory that employs 'is a part of' and 'overlaps.')

3203 Thus, that approach to parsimony is under-motivated. Because of this lack of motivation  
3204 and the aforementioned worries about arbitrariness, many metaphysicians  
3205 reject (NI-Parsimony).

3206 Thus far, I have presented four approaches to parsimony. I have suggested  
3207 that the overall most attractive package for applying parsimony to theory  
3208 choice is one that (i) can include (NO Parsimony), (ii) definitely includes (KO-  
3209 Parsimony) and (KI-Parsimony), but (iii) does not include (NI-Parsimony). Not  
3210 coincidentally, this is a package that has recently gained prominence amongst  
3211 metaphysicians who care about the parsimony of their theories. Even the

---

8 Cf. Cowling (2013); Sider (2011).

9 Some philosophers contest the connection between parsimony and truth—either as it relates to metaphysical theories specifically or as it relates to any descriptive theory. Cf. Brenner (2017); Sober (2015); Willard (2014).

3212 most ardent supporters of parsimony have shied away from including (NI-  
 3213 Parsimony). Sider (2011, 258–259) admits that ‘[t]here is a real question about  
 3214 which of propositional logic’s connectives carve at the joints, and similarly  
 3215 for  $\forall$  and  $\exists$ ,’ and yet nevertheless ‘egalitarian answers can be given. . . [o]ne  
 3216 might hold that both  $\exists$  and  $\forall$  carve at the joints, or that all the truth-functional  
 3217 connectives do, and thus avoid drawing invidious metaphysical distinctions.’<sup>10</sup>

3218 But, as I will now show, there is a puzzle that undermines this package’s  
 3219 credibility.

## 3220 2 The Puzzle

3221 In this section, I will argue that (KO-Parsimony) entails a restricted form of  
 3222 (NI-Parsimony). I will then suggest that this is no mere entailment; properly  
 3223 understood, (KO-Parsimony) and (NI-Parsimony) target the same feature of  
 3224 reality, the structure of the world. Thus, insofar as these two approaches to  
 3225 parsimony are motivated by a desire to posit a simple world, it is puzzling  
 3226 that metaphysicians should treat them so differently.

3227 To illustrate these connections, I will work through a paradigm example  
 3228 of the neo-Quinean methodology at work in the metaphysics of composite  
 3229 objects.

3230 According to compositional nihilism, there are no composite objects—no  
 3231 tables, no chairs, and no people (if people are composite objects). Yet natural  
 3232 language claims like

3233 Some composite objects are larger than other composite objects

3234 seem undeniably true.<sup>11</sup> The most straightforward regimentation of this En-  
 3235 glish claim using first order logic is:

$$\exists x \exists y (C(x) \wedge C(y) \wedge (x \neq y) \wedge L(x, y))$$

10 Here, I avoid using truth-functional operators (like propositional logic’s connectives) as examples, since someone may argue that truth-functional operators aren’t primitive anyway. Instead, truth-functional operators may be defined in terms of their truth-tables, which ultimately depend on primitive notions of truth and falsity. Thanks to an anonymous reviewer for this suggestion.

11 Some metaphysicians (e.g., Merricks 2001) say that such claims are false. Nevertheless, there is a sense in which such claims are “nearly as good as true.” Nothing in what follows depends on the difference between what is true (and later paraphrased) and what is nearly as good as true. Cf. Bennett (2009, 58–59).

3236 which informally reads ‘There is an  $x$  and there is a  $y$  such that  $x$  is a composite  
 3237 object,  $y$  is a composite object,  $x$  is not identical with  $y$ , and  $x$  is larger than  $y$ .’  
 3238 According to orthodox neo-Quineanism, if we endorse this regimentation we  
 3239 thereby incur an ontological commitment to composite objects.

3240 But we want to avoid an ontological commitment to composite objects.  
 3241 This is in part because (KO-Parsimony) recommends reducing the number of  
 3242 posited ontological kinds when feasible. Composite objects form an ontologi-  
 3243 cal kind. So we ought to avoid positing them.

3244 How do we accomplish that goal? It is not enough merely to reduce the num-  
 3245 ber of references to composite objects or to relegate claims about composite  
 3246 objects to a theoretically insignificant role. On the neo-Quinean methodology,  
 3247 we posit an ontological kind when, in stating our theory, we employ a predi-  
 3248 cate that ranges over entities found within that kind. Thus, we need to avoid  
 3249 the mention of composite objects altogether. To accomplish that, we need to  
 3250 find an alternative regimentation to the English sentence (‘Some composite  
 3251 objects are larger than other composite objects.’) that uses only nihilistically  
 3252 acceptable ideology.

3253 Here’s how we can do that. First, we replace the composite object predicate,  
 3254 ‘ $C$ ,’ with the predicate ‘ $AC$ ,’ which reads as ‘arranged composite-object-wise.’  
 3255 This predicate ranges over the things that are spatially distributed as if they  
 3256 composed an object. If contemporary physics is correct, the entities that satisfy  
 3257 this predicate are quarks, leptons, and bosons. But so as to not presuppose  
 3258 any particular theory, let’s call them—whatever they are—“simples.” ‘ $AC$ ’  
 3259 ranges over simples, but in a non-distributive manner. No single simple is  
 3260 arranged composite-object wise. Rather, all of the simples are collectively  
 3261 arranged composite-object-wise. Finally, we must be able to quantify over  
 3262 simples arranged composite-object-wise in a way that avoids committing  
 3263 ourselves to something “over and above” those simples. To that end, we sup-  
 3264 plement first-order logic’s singular quantification with plural quantification.  
 3265 Following some fairly standard notation from Burgess and Rosen (1997), we  
 3266 can use doubled letters (e.g., ‘ $xx$ ,’ ‘ $yy$ ’) to represent the variables for plural  
 3267 quantification. We can then regiment the English sentence as follows:

$$\exists xx \exists yy (AC(xx) \wedge AC(yy) \wedge (xx \neq yy) \wedge L(xx, yy))$$

3268 This sentence successfully avoids an ontological commitment to composite  
 3269 objects.

3270 Yet things are not so simple. We can use plural quantification to eliminate  
 3271 singular references to composite objects. But English also plausibly includes

3272 plural references to composite objects.<sup>12</sup> Consider, for example, the following  
3273 sentence:

3274 Some composite objects are in contact only with one another.

3275 We would need to employ plural quantification in the regimentation of this  
3276 sentence even with an ontological commitment to composite objects. For  
3277 instance, where “*T*” is a predicate that ranges over things in contact and ‘<’ is a  
3278 special relation between individuals and pluralities of individuals, functioning  
3279 like the English expression ‘among’:

$$\begin{aligned} & \exists xx[\forall u((u < xx) \rightarrow C(u)) \wedge \\ & \forall v\forall w(((v < xx) \wedge T(v, w)) \rightarrow ((w < xx) \wedge v \neq w))] \end{aligned}$$

3280 From an ideological perspective, this regimented sentence is already quite ugly.  
3281 But, because it employs a predicate for composite objects, it would commit  
3282 us to the existence of composite objects. So, to avoid such a commitment, we  
3283 must construct a different regimentation that does not use such a predicate.  
3284 This nihilistically acceptable regimentation will be even uglier. That’s because  
3285 it must rely on plurally plural—i.e., perplural—quantification. Just as plural  
3286 quantification ranges over pluralities of individuals, perplural quantification  
3287 ranges over second-level pluralities of pluralities. Let’s use tripled letters (e.g.,  
3288 ‘*xxx*,’ ‘*yyy*’) to represent the variables for perplural quantification. We then  
3289 get the following regimentation:

$$\begin{aligned} & \exists xxx[\forall uu((uu < xxx) \rightarrow AC(uu)) \wedge \\ & \forall v\forall w(((v < xxx) \wedge T(vv, ww)) \rightarrow ((w < xxx) \wedge vv \neq ww))] \end{aligned}$$

3290 In this way, metaphysicians can avoid an ontological commitment to com-  
3291 posite objects, thereby minimizing the kinds of objects to which they are  
3292 ontologically committed. But their use of (primitive) perplural quantification  
3293 increases the ideological kinds to which they are committed.

3294 So far as the metaphysics of composite objects goes, we have two options.  
3295 First, we can employ a predicate that ranges over composite objects. Or, to  
3296 avoid the ontological commitment, we can remove the predicate. Choosing  
3297 this second option seems to involve a trade-off between a specially problematic  
3298 predicate and a more complicated form of quantification.

12 It is contentious whether English contains genuine perplural locutions (see Linnebo and Nicolas 2008; McKay 2006, 46–52). I cannot speak to other natural languages.

3299 Our intuitive grasp of the relevant concepts initially suggested that ontology  
3300 and ideology are quite distinct. So it's surprising that a commitment to (KO-  
3301 Parsimony) entails a *de facto* commitment to (NI-Parsimony). This connection  
3302 cries out for explanation.

3303 In fact, I think the explanation is quite straightforward for ideological  
3304 externalists. If we use a theory's ideology to pick out features of the world,  
3305 then it's entirely plausible that in doing so we sometimes pick out ontological  
3306 kinds.

3307 Think of it this way. The elimination of a single object from a metaphysi-  
3308 cian's ontology improves its quantitative ontological parsimony. So, too, does  
3309 the elimination of every object of a given kind. But the elimination of an  
3310 ontological kind does not necessarily result in the elimination of any ob-  
3311 jects. It's perfectly ordinary for a reductive project to "relocate" the objects  
3312 of one kind into the province of another. For example, David Lewis's modal  
3313 realism (1986) is ontologically parsimonious insofar as it avoids an ontolog-  
3314 ical commitment to *sui generis* possible worlds. But it does not minimize  
3315 the overall number of objects; in a manner of speaking, what would have  
3316 been *sui generis* possible worlds are instead causally isolated concrete enti-  
3317 ties. So, (KO-Parsimony) should not be understood as an efficient means of  
3318 reducing the overall number of objects posited. Similarly, (KO-Parsimony)  
3319 should not be understood merely as a preference for "empty kinds" over "pop-  
3320 ulated kinds." In many cases, whether or not an ontological kind is populated  
3321 should depend on contingent facts of the world rather than metaphysical  
3322 necessities. (KO-Parsimony) should be understood as a preference for the *elim-*  
3323 *ination* of ontological kinds. As the compositional example above suggests,  
3324 the elimination of an ontological kind is achieved by the abandonment of its  
3325 corresponding predicate. Here is where ideological externalism is relevant.  
3326 When a theory commits to an ontological kind, it is not committing to some  
3327 object that it quantifies over. Rather, when a theory commits to an ontological  
3328 kind, it is committing to a structural feature of the world that corresponds  
3329 to a predicate employed by the theory's ideology. Similarly, when a theory  
3330 eliminates an ontological kind, it eliminates a structural feature of the world.  
3331 Ontological kinds are features of the world's metaphysical structure.

3332 Compare this theoretical identification to the theoretical identification of  
3333 water and H<sub>2</sub>O. Our concept of water is quite different from our concept of  
3334 H<sub>2</sub>O: our concept of water predates our concept of H<sub>2</sub>O; our concept of water  
3335 is rooted in its geographic, biological, and sociological functions whereas our  
3336 concept of H<sub>2</sub>O is rooted in the scientific discipline of chemistry; and so on.

3337 As a matter of fact, though, the two concepts pick out the same substance. Of  
 3338 course, in some sense our concept of water “could have” picked out a different  
 3339 substance. Perhaps, even, our concept of water “could have” picked out a  
 3340 metaphysically gruesome disjunction of substances. But that’s not how things  
 3341 turned out. Consequently, to be concerned with water is to be concerned with  
 3342  $H_2O$ . Imagine someone who stressed the importance of bringing water on a  
 3343 camping trip. If they stressed the importance of bringing *water* but denied  
 3344 the importance of bringing  $H_2O$ , we would be confused—and rightly so.

3345 So, too, for ontological kinds and the world’s metaphysical structure. While  
 3346 our concept of an ontological kind may predate our concept of the world’s  
 3347 metaphysical structure, the two concepts ultimately pick out the same fea-  
 3348 ture. Of course, there may be some differences between the two theoretical  
 3349 identifications. Those who maintain a firm distinction between the *a priori*  
 3350 and the *a posteriori* would likely consider “Water is  $H_2O$ ” to be an *a posteriori*  
 3351 identification and “Ontological kinds are metaphysical structure” to be an  
 3352 *a priori* identification. But, assuming the identities hold, many of the com-  
 3353 parisons are apt. If a metaphysician stresses the importance of minimizing  
 3354 the ontological kinds posited by a theory, we should expect them to stress the  
 3355 importance of minimizing the structural complexity posited by a theory—it’s  
 3356 the same thing that is being minimized! At a minimum, the metaphysician  
 3357 owes us an explanation for the difference in attitude.

3358 Thus far, I have argued that those committed to (KO-Parsimony) should be  
 3359 committed to a restricted version of (NI-Parsimony). I have also suggested that  
 3360 there is an identity between the targets of these two principles of parsimony;  
 3361 both seek to minimize the structural complexity of the world. It does not  
 3362 follow that qualitative ontological parsimony *just is* quantitative ideological  
 3363 parsimony. There will still be instances of the latter that aren’t instances of  
 3364 the former. Consider, for instance, a choice between two competing modal  
 3365 theories. Some forms of actualism (like those in [Prior and Fine 1977](#)) eschew  
 3366 quantifying over possible worlds and take the sentential modal operators as  
 3367 primitive. Suppose that actualist theory  $T_1$  takes both ‘ $\Box$ ’ and ‘ $\Diamond$ ’ as primitive  
 3368 and actualist theory  $T_2$  takes only ‘ $\Box$ ’ as primitive, defining ‘ $\Diamond$ ’ in the standard  
 3369 way. (NI-Parsimony) would recommend  $T_1$  over  $T_2$  because it employs one  
 3370 less bit of ideology. But by hypothesis neither theory posits more or fewer  
 3371 kinds of entities. Thus, some disputes about ideology are not reducible to  
 3372 disputes that involve ontology.<sup>13</sup>

---

13 I develop this point more fully in Finocchiaro (2019b).

3373 Here's a small, but important, complication that I've ignored.<sup>14</sup> Thus far, I  
3374 have worked through a single case, the metaphysics of composition. Even if  
3375 what I have said holds for this case, does the point generalize? Or is it merely  
3376 an artifact of the case that might or might not apply to others?

3377 The point generalizes. On the neo-Quinean paradigm, there is no ontologi-  
3378 cal commitment to something unless there is a regimented sentence held to be  
3379 true which includes a bound variable that must refer to that thing. But there is  
3380 no need to have such a referring bound variable unless that variable attaches  
3381 to a predicate of some kind. In other words, because ontological parsimony is  
3382 a difference in ontology and because ontology is always expressed through  
3383 ideology, ontological parsimony always involves a difference in ideology.

3384 There is one slight exception. Some metaphysicians adopt principles of par-  
3385 simony that discriminate on the basis of fundamentality. For example, Schaffer  
3386 (2009) adopts the Laser, which recommends minimizing the number of *fun-*  
3387 *damental* entities but does not care about the number of *non-fundamental*  
3388 entities. Such a principle makes the connection between ontology and ideology  
3389 weaker. More specifically, when using the Laser there will be predicates—  
3390 the ones corresponding to non-fundamental entities—whose elimination or  
3391 introduction would not impact ontological parsimony.

3392 But this exception does not solve the puzzle. First, it's unclear what the  
3393 status of such predicates is. Plausibly, non-fundamental ontology is expressed  
3394 through non-fundamental ideology. If so, then this exception is simply irrele-  
3395 vant to the puzzle I've presented. Second, this exception still entails a strong  
3396 relationship between fundamental ontology and fundamental ideology. So, at  
3397 best, it would solve only part of the puzzle.

### 3398 **3 What to Do?**

3399 I will end by briefly discussing three ways to respond to the puzzle about  
3400 parsimony. Each has its advantages and disadvantages. While I do favor one  
3401 of the ways over the others, I think all three are worth developing more fully.

3402 First, we could try to resist the puzzle. I generated the puzzle by assum-  
3403 ing orthodox neo-Quineanism. One way of resisting, then, is to reject the  
3404 claim that a theory's ontology is that over which the theory quantifies. There  
3405 are several alternatives to the Quinean criterion of ontological commitment,  
3406 but one promising option is the truthmaker view. On the truthmaker view,

---

14 Thanks to an anonymous reviewer for pushing me to address this issue.

3407 a theory's ontology is that which makes the theory's sentences true.<sup>15</sup> Im-  
 3408 portantly, the view explicitly permits two theories to differ with respect to  
 3409 their ideologies without also differing with respect to their ontological com-  
 3410 mitments. For instance, on the truthmaker view a theory might truly state  
 3411 "Some composite objects are larger than other composite objects" without  
 3412 incurring an ontological commitment to composite objects. What matters is  
 3413 not what the sentence quantifies over but rather what makes the sentence  
 3414 true—and what makes the sentence true need not be composite objects. More  
 3415 importantly, the view entails that the two regimentations offered above—  
 3416 " $\exists x \exists y (C(x) \wedge C(y) \wedge (x \neq y) \wedge L(x, y))$ " and " $\exists x x \exists y y (AC(xx) \wedge AC(yy) \wedge$   
 3417  $(xx \neq yy) \wedge L(xx, yy))$ "—have the same ontological commitments. The  
 3418 change in ideology does not impact the ontology. Thus, on the truthmaker  
 3419 view of ontological commitment, (KO-Parsimony) does not entail any version  
 3420 of (NI-Parsimony), nor does it suggest an identity between their targets. In a  
 3421 way, then, the puzzle about parsimony could motivate us to reject orthodox  
 3422 neo-Quineanism.

3423 Those of us not yet ready to abandon orthodoxy have to either embrace  
 3424 the puzzle or downplay its significance. I suspect many would prefer the  
 3425 second option. Some metaphysicians (e.g., Bennett 2009) have characterized  
 3426 many metaphysical disputes as being, as bottom, trade-offs between ontology  
 3427 and ideology. This characterization is hard to maintain if they have the same  
 3428 target (i.e., the world's metaphysical structure). It seems, then, that my puzzle  
 3429 puts that characterization in a hard place. But perhaps the essence of their  
 3430 characterization can be maintained. I can see two strategies for doing so.

3431 On the first strategy, there are many more ideological kinds than previously  
 3432 assumed. More specifically, each predicate that expresses an ontological kind  
 3433 forms its own ideological kind. If this is so, then (KO-Parsimony) actually  
 3434 entails (KI-Parsimony) and the methodological tension vanishes. But here's a  
 3435 challenge that this strategy must overcome. By following the neo-Quinean  
 3436 orthodoxy, we eliminate ontologically committing predicates but we do not  
 3437 eliminate the complements of those predicates. So, for instance, the com-  
 3438 positional nihilist eliminates 'composite object' but does not eliminate 'not  
 3439 a composite object,' otherwise known as 'simple.' Yet, intuitively, "positive"

---

15 See Rettler (2016, 21). Rettler even appears to gesture toward a version of my puzzle when he says, "[I]t's true, just looking at the sentences will no longer tell you which theory wins the day with respect to parsimony of ontological commitments. But it never should have." In what follows I will simplify my discussion by ignoring Rettler's distinction between the general truthmaker view and the specific truthmaker view.



3440 predicates like ‘composite object’ and “negative” complements like ‘simple’ are  
3441 of the same ideological kind. So, those who want to pursue this first strategy  
3442 of downplaying the significance of the puzzle must offer a more sophisticated  
3443 means of individuating ideological kinds.

3444 On the second strategy, there are two categories of ideology such that (i)  
3445 we ought to minimize the number of ideological bits from the first category,  
3446 and (ii) it is not the case that we ought to minimize the number of ideological  
3447 bits from the second category. Obviously, those who pursue this strategy must  
3448 offer some explanation for the difference in treatment. One somewhat radical  
3449 explanation is to say that structural simplicity is more important in some  
3450 domains than it is in others. I do not see how this explanation can be plausibly  
3451 maintained. Parsimony is currently treated as a comprehensive value: choose  
3452 the theory that is *overall* more simple. Why would simplicity in one domain be  
3453 less important (i.e., less truth conducive) than simplicity in another domain?  
3454 On an alternative explanation, the relationship between ideological bits and  
3455 metaphysical structure is more nuanced than previously thought. Perhaps  
3456 ideological bits are more fine-grained than the corresponding structure. If so,  
3457 then some ideological bits (like ‘is a part of’ and ‘overlaps’) would correspond  
3458 to the same aspect of the world’s metaphysical structure, and so there is  
3459 no need to choose between the two. In contrast, other ideological bits (like  
3460 ‘composite object’ and ‘simple’) would correspond to different aspects of the  
3461 world’s metaphysical structure, and so there is value in eliminating one if not  
3462 the other. This explanation is interesting. But as it stands it is *ad hoc*. In the  
3463 absence of a worked-out account of ideological correspondence, why should  
3464 we think that it works the way this strategy needs it to work?

3465 That leaves the third response: embrace the puzzle. If we embrace the  
3466 puzzle, we ought to claim that (NI-Parsimony) is no less justified a principle  
3467 than (KO-Parsimony). This claim is quite shocking (well, as shocking as an  
3468 esoteric claim about the proper methodological application of parsimony  
3469 can be, anyway). (KO-Parsimony) has a rich history and is likely the most  
3470 broadly endorsed approach to parsimony. In contrast, almost no one explicitly  
3471 endorses (NI-Parsimony). Nevertheless, by embracing the puzzle we can save  
3472 neo-Quineanism as well as the standard characterization of metaphysical  
3473 disputes as disputes that involve trade-offs between ontology and ideology.  
3474 Yet those who pursue this third strategy have their own explaining to do.  
3475 Intuitively, it seems objectionably arbitrary to choose between functionally  
3476 equivalent terminology. So why isn’t it? For example, why should we reduce  
3477 the number of compositional predicates we employ in stating our theories?

3478 Perhaps we can extend the standard motivations for parsimony-based  
 3479 considerations and say that we should reduce the number of compositional pred-  
 3480 icates because the resulting theory posits a more simple structure and is  
 3481 therefore more likely to be true. This might still generate an epistemic dead-  
 3482 lock with regard to competing “equivalent” theories. (NI-Parsimony) would  
 3483 suggest that a theory that employs only ‘overlaps’ is more likely to accurately  
 3484 represent the compositional structure of the world than a theory that employs  
 3485 both ‘overlaps’ and ‘is a part of.’ *Mutatis mutandis* for a theory that employs  
 3486 only ‘is a part of.’ But at this point our methodology fails us and we do not  
 3487 know which of the two predicates we ought to employ.<sup>16</sup>

3488 Personally, I think we ought to embrace the puzzle. It’s not a perfect re-  
 3489 sponse, but it is the best available. Neo-Quineanism is battle-tested orthodoxy.  
 3490 (More modestly, neo-Quineanism is much closer to the center of my web  
 3491 of belief than are the other elements of the puzzle.) For that reason I reject  
 3492 the first response. The second response raises a number of issues regarding  
 3493 ideological correspondence. I am doubtful that those issues can be addressed  
 3494 satisfactorily. So I also reject the second response. Finally, I do not think that  
 3495 the third response is that bad. I don’t know how to choose between overlap  
 3496 and parthood. I don’t even know how to think about that choice. But a hard  
 3497 choice is not *ipso facto* a bad choice.\*

3498 Peter Finocchiaro  
 3499  0000-0003-4060-7061  
 3500 Wuhan University  
 3501 peter.w.finocchiaro@gmail.com

## 3502 References

- 3503 BENNETT, Karen. 2009. “Composition, Colocation, and Metaontology.” in *Metameta-*  
 3504 *physics: New Essays on the Foundations of Ontology*, edited by David J. CHALMERS,  
 3505 David MANLEY, and Ryan WASSERMAN, pp. 38–76. Oxford: Oxford University  
 3506 Press, doi:10.1093/oso/9780199546046.003.0002.
- 3507 BRENNER, Andrew. 2017. “Simplicity as a Criterion of Theory Choice in Metaphysics.”  
 3508 *Philosophical Studies* 174(11): 2687–2707, doi:10.1007/s11098-016-0805-1.

16 Cf. McSweeney (2019, 127–128).

\* For their help in the development of this paper, I thank Andrew Brenner, Rebecca Chan, Alessandro Torza, Jeffrey J. Watson, Qiong Wu, the audience at my 2020 Central APA conference session, and two anonymous reviewers. I also thank Collegetown Bagels for their bagels, coffee, and outdoor seating, all of which are indispensable for writing during a pandemic.

- 3509 BURGESS, John P. and ROSEN, Gideon. 1997. *A Subject with No Object: Strategies*  
 3510 *for Nominalistic Interpretation of Mathematics*. Oxford: Oxford University Press,  
 3511 doi:10.1093/0198250126.001.0001.
- 3512 CAMERON, Ross P. 2012. "Why Lewis's Analysis of Modality Succeeds in Its Reductive  
 3513 Ambitions." *Philosophers' Imprint* 12(8), [http://hdl.handle.net/2027/spo.3521354.](http://hdl.handle.net/2027/spo.3521354.0012.008)  
 3514 0012.008.
- 3515 CHALMERS, David J., MANLEY, David and WASSERMAN, Ryan, eds. 2009. *Metameta-*  
 3516 *physics: New Essays on the Foundations of Ontology*. Oxford: Oxford University  
 3517 Press, doi:10.1093/oso/9780199546046.001.0001.
- 3518 CHURCHLAND, Paul M. 1984. *Matter and Consciousness*. Cambridge, Massachusetts:  
 3519 The MIT Press.
- 3520 COWLING, Sam. 2013. "Ideological Parsimony." *Synthese* 190(17): 3889–3908, doi:10.1  
 3521 007/s11229-012-0231-7.
- 3522 FINOCCHIARO, Peter. 2019a. "The Explosion of Being: Ideological Kinds in Theory  
 3523 Choice." *The Philosophical Quarterly* 69(276): 486–510, doi:10.1093/pq/pqz005.
- 3524 —. 2019b. "The Intelligibility of Metaphysical Structure." *Philosophical Studies* 176(3):  
 3525 581–606, doi:10.1007/s11098-017-1029-8.
- 3526 —. 2021. "Ideology and Its Role in Metaphysics." *Synthese* 198(2): 957–983, doi:10.100  
 3527 7/s11229-018-02077-6.
- 3528 HORGAN, Terence E. and POTRČ, Matjaž. 2008. *Austere Realism: Contextual Semantics*  
 3529 *Meets Minimal Ontology*. Cambridge, Massachusetts: The MIT Press, doi:10.7551/  
 3530 mitpress/9780262083768.001.0001.
- 3531 VAN INWAGEN, Peter. 1998. "Meta-Ontology." *Erkenntnis* 48(2/3): 233–250. Reprinted  
 3532 in van Inwagen (2001, 13–31), doi:10.1023/A:1005323618026.
- 3533 —. 2001. *Ontology, Identity, and Modality: Essays in Metaphysics*. Cambridge: Cam-  
 3534 bridge University Press.
- 3535 LEWIS, David. 1973. *Counterfactuals*. Cambridge, Massachusetts: Harvard University  
 3536 Press. Cited after republication as Lewis (2001).
- 3537 —. 1983. *Philosophical Papers, Volume 1*. Oxford: Oxford University Press, doi:10.109  
 3538 3/0195032047.001.0001.
- 3539 —. 1986. *On the Plurality of Worlds*. Oxford: Blackwell Publishers.
- 3540 —. 2001. *Counterfactuals*. Oxford: Blackwell Publishers. Republication of Lewis (1973).
- 3541 LEWIS, David and LEWIS, Stephanie R. 1970. "Holes." *Australasian Journal of Philoso-*  
 3542 *phy* 48(2): 206–212. Reprinted in Lewis (1983, 3–9), doi:10.1080/00048407012341  
 3543 181.
- 3544 LINNEBO, Øystein and NICOLAS, David. 2008. "Superplurals in English." *Analysis*  
 3545 68(3): 186–197, doi:10.1093/analys/68.3.186.
- 3546 MCKAY, Thomas J. 2006. *Plural Predication*. Oxford: Oxford University Press, doi:10.1  
 3547 093/acprof:oso/9780199278145.001.0001.
- 3548 MCSWEENEY, Michaela Markham. 2019. "Following Logical Realism where it Leads."  
 3549 *Philosophical Studies* 176(1): 117–139, doi:10.1007/s11098-017-1008-0.

- 3550 MELIA, Joseph. 2008. "A World of Concrete Particulars." in *Oxford Studies in Meta-*  
 3551 *physics*, volume IV, edited by Dean W. ZIMMERMAN, pp. 99–124. Oxford: Oxford  
 3552 University Press, doi:10.1093/oso/9780199542987.003.0005.
- 3553 MERRICKS, Trenton. 2001. *Objects and Persons*. Oxford: Oxford University Press, doi:10  
 3554 .1093/0199245363.001.0001.
- 3555 NOLAN, Daniel Patrick. 1997. "Quantitative Parsimony." *The British Journal for the*  
 3556 *Philosophy of Science* 48(3): 329–343, doi:10.1093/bjps/48.3.329.
- 3557 PARSONS, Josh. 2014. "The Many Primitives of Mereology." in *Mereology and Location*,  
 3558 edited by Shieva KLEINSCHMIDT, pp. 3–12. Oxford: Oxford University Press, doi:10  
 3559 .1093/acprof:oso/9780199593828.003.0001.
- 3560 PAUL, Laurie A. 2017. "A One Category Ontology." in *Being, Freedom, and Method:*  
 3561 *Themes from the Philosophy of Peter van Inwagen*, edited by John Adorno KELLER,  
 3562 pp. 32–61. Oxford: Oxford University Press, doi:10.1093/acprof:oso/978019871570  
 3563 2.003.0003.
- 3564 PRIOR, Arthur Norman and FINE, Kit. 1977. *Worlds, Times and Selves*. London: Gerald  
 3565 Duckworth & Co. Edited by Kit Fine; based on manuscripts by Prior with a preface  
 3566 and a postscript by Kit Fine.
- 3567 QUINE, Willard van Orman. 1948. "On What There Is." *The Review of Metaphysics* 2(5):  
 3568 21–38. Republished as Quine (1951), reprinted in Quine (1953, 1–19).
- 3569 —. 1951. "On What There Is." *Proceedings of the Aristotelian Society, Supplementary*  
 3570 *Volume* 25: 217–234. Reprint of Quine (1948), doi:10.1093/aristoteliansupp/25.1.  
 3571 125.
- 3572 —. 1953. *From a Logical Point of View: 9 Logico-Philosophical Essays*. Cambridge,  
 3573 Massachusetts: Harvard University Press. Cited after the revised edition: Quine  
 3574 (1961).
- 3575 —. 1961. *From a Logical Point of View: 9 Logico-Philosophical Essays*. 2nd ed. Cam-  
 3576 bridge, Massachusetts: Harvard University Press. Revised edition of Quine (1953),  
 3577 reprinted 1980.
- 3578 RETTLER, Bradley. 2016. "The General Truthmaker View of Ontological Commitment."  
 3579 *Philosophical Studies* 173(5): 1405–1425, 1427, doi:10.1007/s11098-015-0526-x.
- 3580 SCHAFFER, Jonathan. 2009. "On What Grounds What." in *Metametaphysics: New Essays*  
 3581 *on the Foundations of Ontology*, edited by David J. CHALMERS, David MANLEY,  
 3582 and Ryan WASSERMAN, pp. 347–383. Oxford: Oxford University Press, doi:10.109  
 3583 3/oso/9780199546046.003.0012.
- 3584 —. 2015. "What Not to Multiply Without Necessity." *Australasian Journal of Philosophy*  
 3585 93(4): 644–664, doi:10.1080/00048402.2014.992447.
- 3586 SIDER, Theodore. 2011. *Writing the Book of the World*. Oxford: Oxford University Press,  
 3587 doi:10.1093/acprof:oso/9780199697908.001.0001.
- 3588 —. 2013. "Against Parthood." in *Oxford Studies in Metaphysics*, volume VIII, edited  
 3589 by Karen BENNETT and Dean W. ZIMMERMAN, pp. 237–293. New York: Oxford  
 3590 University Press, doi:10.1093/acprof:oso/9780199682904.003.0006.

- 3591 SOBER, Elliott R. 2015. *Ockham's Razors: A User's Manual*. Cambridge: Cambridge  
3592 University Press, doi:10.1017/cbo9781107705937.
- 3593 TALLANT, Jonathan. 2013. "Quantitative Parsimony and the Metaphysics of Time:  
3594 Motivating Presentism." *Philosophy and Phenomenological Research* 87(3): 688–795,  
3595 doi:10.1111/j.1933-1592.2012.00617.x.
- 3596 TURNER, Jason. 2015. "What's so Bad About Second-Order Logic?" in *Quantifiers,*  
3597 *Quantifiers, and Quantifiers: Themes in Logic, Metaphysics and Language*, edited  
3598 by Alessandro TORZA, pp. 463–488. Synthese Library n. 373. Dordrecht: Springer,  
3599 doi:10.1007/978-3-319-18362-6\_20.
- 3600 WILLARD, Mary Beth. 2014. "Against Simplicity." *Philosophical Studies* 167(1): 165–181,  
3601 doi:10.1007/s11098-013-0228-1.

PROOF

PROOF

3602 David Armstrong on the Metaphysics  
3603 of Mathematics

TOM DONALDSON

3604 This paper has two components. The first, longer component (sections  
3605 1–6) is a critical exposition of Armstrong’s views about the metaphysics of  
3606 mathematics, as they are presented in *Truth and Truthmakers* and *Sketch*  
3607 *for a Systematic Metaphysics*. In particular, I discuss Armstrong’s views  
3608 about the nature of the cardinal numbers, and his account of how modal  
3609 truths are made true. In the second component of the paper (section 7),  
3610 which is shorter and more tentative, I sketch an alternative account of  
3611 the metaphysics of mathematics. I suggest we insist that mathematical  
3612 truths have physical truthmakers, without insisting that mathematical  
3613 objects themselves are part of the physical world.

3614 A prime number  $p$  is a “Sophie Germain prime” if  $(2p + 1)$  is also prime.  
3615 It is conjectured that there exist infinitely many Sophie Germain primes.  
3616 I don’t know whether this conjecture is true, but what I do know is that  
3617 there exist *some* Sophie Germain primes: 2 is an example; 3 is another;  
3618  $(2, 618, 163, 402, 417 \times 2^{1290000} - 1)$  is a third, or so I am told. Now it is  
3619 obvious that every Sophie Germain prime is a number; and it follows that  
3620 there exist some numbers—or so it seems.

3621 But what kind of a thing is a number?

3622 This is a difficult question, but one point at least seems clear: numbers  
3623 and other mathematical entities are “abstract,” in the sense that they have no  
3624 causal powers and no location in spacetime. We are told that the number zero  
3625 was discovered in India, but it would be a mistake to go to India *now* to look  
3626 for it—and not because it has subsequently been moved. You can’t trip over  
3627 the number three. The polynomial  $(x^2 - 3x + 2)$  can be split into two factors,  
3628  $(x - 2)$  and  $(x - 1)$ , but not by firing integers at it in a particle accelerator. The  
3629 empty set has no gravitational field. And so on.

3630 And so, to labour the point, it seems that some abstract entities exist.

3631 And yet David Armstrong began his last book by endorsing what he called  
3632 “naturalism”:

3633 I begin with the assumption that all that exists is the space-time  
3634 world, the physical world as we say. [...] [This] means the rejection  
3635 of what many contemporary philosophers call “abstract objects”,  
3636 meaning such things as numbers or Platonic Forms or classes,  
3637 where these are supposed to exist “outside of” or “extra to” space-  
3638 time. (2010, 1)

3639 Despite this naturalism, Armstrong did not reject any part of mainstream  
3640 mathematics. Indeed, he insisted that the truths of orthodox pure mathematics  
3641 are necessary and *a priori* (2010, ch. 12).

3642 In this paper, I explore Armstrong’s attempt to reconcile his denial of ab-  
3643 stract entities with his commitment to orthodox mathematics. To be more  
3644 specific, the paper has three goals.

- 3645 1. Armstrong wrote a vast amount on mathematics, and this writing is  
3646 spread among many papers and books. Some of this work is complex,  
3647 and Armstrong changed his mind on certain important questions. My  
3648 first goal in this paper is to describe—clearly, briefly, and in one place—  
3649 Armstrong’s mature views on the metaphysics of mathematics, includ-  
3650 ing relevant aspects of his work on the metaphysics of modality. To  
3651 prevent the discussion from sprawling, I focus particularly on Arm-  
3652 strong’s account of cardinal number, as it is presented in his last two  
3653 books: *Truth and Truthmakers* (2004, “T&T”) and *Sketch for a Systematic*  
3654 *Metaphysics* (2010, “SSM”). Section 1 and section 2 describe Armstrong’s  
3655 views about the cardinal numbers; sections 3–6 focus on modality.
- 3656 2. My second goal is to present some novel—and, I believe, definitive—  
3657 objections to Armstrong’s views on the metaphysics of mathematics.  
3658 These objections are presented in sections 5 and 6.
- 3659 3. My third goal is to recommend a different way of thinking about the  
3660 metaphysics of mathematics—an approach which will, I hope, appeal  
3661 to people who admire Armstrong’s work. Briefly, I will suggest that we  
3662 insist that every mathematical truth has a truthmaker in the physical  
3663 world, without also insisting that mathematical objects themselves are  
3664 physical things. The proposal is presented in more detail in section 7.



## 3665 **1 Cardinal Numbers as Concrete Entities**

3666 The claim that numbers have no spatial location is familiar to metaphysicians,  
 3667 but it sometimes comes as a surprise to students. “But there are three pens  
 3668 on my desk *right now!*” they exclaim, implying that the number three itself is  
 3669 within arm’s reach. According to Armstrong, the surprised students are on to  
 3670 something.

3671 While he rejected “Platonic forms” which exist outside spacetime, Arm-  
 3672 strong did believe that there are properties which exist *within* the particulars  
 3673 that instantiate them (SSM, ch. 2). For Armstrong, there exists a property *is*  
 3674 *red* which exists within London buses, ripe tomatoes, and male cardinals. As  
 3675 he sometimes put it, properties are “immanent” rather than “transcendent.”  
 3676 He inferred that all properties are instantiated. Uninstantiated properties have  
 3677 no place in spacetime, and so no place in Armstrong’s philosophical system  
 3678 (SSM, 15–16). He made the same claim about relations (SSM, 23).

3679 For Armstrong, a cardinal number is a relation between a particular and  
 3680 a property. Specifically, the cardinal number  $\kappa$  is a relation that a particular  
 3681  $x$  bears to a property  $P$  just in case  $x$  has, as mereological parts, exactly  $\kappa$   
 3682 particulars which instantiate  $P$ . A normal octopus bears the *one* relation to  
 3683 the property *is an octopus*, and the *eight* relation to the property *is a limb*.  
 3684 The mereological sum of two normal octopuses bears the *two* relation to the  
 3685 property *is an octopus* and the *sixteen* relation to the property *is a limb*. And  
 3686 so on.<sup>1, 2</sup>

3687 On Armstrong’s view, then, the surprised student is correct to think that  
 3688 cardinal numbers exist within the “spacetime world.” It turns out, then, that  
 3689 one can coherently maintain that numbers exist and that there are no abstract  
 3690 entities.

3691 This is a striking result—and yet, a problem looms. As I said, Armstrong  
 3692 insisted that properties and relations exist only when they are instantiated.  
 3693 On this view, the number  $10^{10^{10}}$  exists only if the spacetime world happens to

1 Armstrong’s theory of cardinal number is closely related to that presented in Kessler (1980). Simons (1982) raises some important objections to Kessler’s account. I believe that Armstrong’s theory is not vulnerable to Kessler’s objections, though I will not pursue the issue here.

2 Note that the property  $P$  will in many cases be a “second-rate property” rather than a genuine universal (SSM, 19). For example, there are seven medium-sized red spoons in my apartment. As Armstrong would put it, the fusion of the contents of my apartment bears the *seven* relation to the property *medium-sized red spoon*. But Armstrong would surely deny that *medium-sized red spoon* is a universal. Armstrong understood that we need “second-rate” properties in cases like this (T&T, 72).

3694 contain  $10^{10^{10}}$  particulars. And it is far from clear that the spacetime world is  
 3695 that large. I call this “the problem of size.”

3696 It is tempting to reply to this objection by insisting that space is infinitely  
 3697 divisible. Discussing the problem of size as it arises for Aristotle, Jonathan  
 3698 Barnes writes:

3699 Physical objects are, in Aristotle’s view, infinitely divisible. That  
 3700 fact ensures that, even within the actual finite universe, we shall  
 3701 always be able to find a group of  $k$  objects, for any  $k$  [...] If the  
 3702 universe consisted simply of a single sphere, it would also contain  
 3703 two objects (two hemispheres), three objects (three third-spheres)  
 3704 and so on. We shall never run short of numbers of things [...].  
 3705 (1985, 122)

3706 This cannot be considered a satisfactory solution to *Armstrong’s* problem,  
 3707 however. For one thing, it is far from clear that Aristotle was correct in thinking  
 3708 that space is infinitely divisible: those who study quantum gravity have been  
 3709 known to speculate that space is in fact discrete. And so the proposed solution  
 3710 is somewhat “hostage to fortune.” More importantly, even if the proposed  
 3711 account does secure the existence of large finite numbers such as  $10^{10^{10}}$ ,  
 3712 it still leaves the existence of transfinite cardinals open to doubt. Standard  
 3713 mathematical descriptions of spacetime (which do imply infinite divisibility)  
 3714 entail that the set of spacetime points has cardinality  $\aleph_1$ . Such accounts leave  
 3715 it unclear whether larger cardinal numbers (e.g.,  $\aleph_2$ ,  $\aleph_3$ , or even  $\aleph_\omega$ ) are  
 3716 instantiated in the physical world—and yet these larger cardinals are very  
 3717 much a part of orthodox mathematics.<sup>3</sup>

3718 In passing, I note that *Armstrong* faced a problem of size in his account of  
 3719 set theory too. While *Armstrong* identified cardinal numbers with relations,  
 3720 he identified sets with individuals. The central claim in *Armstrong’s* account  
 3721 was David Lewis’s “brilliant insight” (T&T, 120): the mereological parts of a  
 3722 set are precisely its non-empty subsets. For example, Lewis’s claim implies  
 3723 that the mereological proper parts of  $\{a, b, c\}$  are  $\{a\}$ ,  $\{b\}$ ,  $\{c\}$ ,  $\{a, b\}$ ,  $\{b, c\}$ , and  
 3724  $\{c, a\}$ . This implies that every singleton set is mereologically simple. Now  
 3725 the set theorists tell us that for each cardinal number  $\kappa$ , there are at least  $\kappa$   
 3726 singletons. Thus, *Armstrong* is stuck with the claim that, within the spacetime  
 3727 world, there are least  $\kappa$  mereologically simple individuals, for each  $\kappa$ . And this  
 3728 seems highly doubtful. Perhaps one could plausibly argue that there are  $\aleph_1$

3 For an introduction to the mathematics of these enormous numbers, see Yarnelle (1964).

3729 mereological simples by saying that each spacetime point is mereologically  
 3730 simple and there are  $\aleph_1$  of those. But it is hard to see any justification for  
 3731 the claim that there are  $\aleph_2$ , or  $\aleph_3$ , or even  $\aleph_\omega$  mereological simples in the  
 3732 spacetime world. Once again, the physical universe seems to be too small to  
 3733 accommodate the ontology of mathematics.<sup>4</sup>

## 3732 2 Armstrong's Possibilism

3735 Armstrong was aware of the problem of size. He responded by claiming that  
 3736 larger cardinal numbers exist *in posse* though not *in esse*:

3737 A Platonist will solve the problem [of size] by postulating unin-  
 3738 stantiated numbers. [...] [However,] my own hard choice [...] is  
 3739 to accept a deflationary doctrine of what it is for a mathemati-  
 3740 cal entity to exist. Plenty of mathematical structures exist in a  
 3741 straightforward sense, because they are instantiated. We can call  
 3742 them *empirical* mathematical structures. [...] But mathematical  
 3743 existence itself, I suggest, should be reckoned as something less.  
 3744 A mathematical entity exists if and only if it is *possible* that there  
 3745 should be instantiations of that structure. (T&T, 117)

3746 Armstrong called this *possibilism*. The doctrine is not easy to interpret. In the  
 3747 above passage, Armstrong seems to suggest that a mathematical entity exists  
 3748 provided that it *could* be instantiated—even if it is in fact not instantiated.  
 3749 However, in this same passage, Armstrong *contrasts* his possibilism with the  
 3750 “Platonist” claim that there are uninstantiated mathematical entities.

3751 The following quotation gives us a clue about what Armstrong meant:

3752 We say ‘ $7 + 5 = 12$ ’, but this can be rendered more transparently,  
 3753 though more boringly, as (necessarily, if there are seven things  
 3754 and five further things, then the sum of these things are twelve  
 3755 things). (T&T, 101)

3756 What this passage suggests is that, for Armstrong, while the sentence “ $7 + 5 =$   
 3757  $12$ ” appears to describe a relation among three mathematical entities (*viz.* the  
 3758 numbers seven, five and twelve) it is in fact a generalization about pluralities  
 3759 of marbles, pebbles, sticks, or whatever. More generally, the possibilist main-  
 3760 tains that while pure mathematics appears to describe a domain of special

---

4 For a more thorough discussion of this point, see Rosen (1995).

3761 mathematical entities, it in fact consists of *modal* statements—statements  
3762 about what is necessary or possible.

3763 Now Armstrong did not develop this proposal systematically; instead, he  
3764 endorsed Geoffrey Hellman’s modal structuralism:

3765 I recognize, of course, that asserting here this [...] doctrine of  
3766 mathematical existence is to a degree a matter of hand-waving. I  
3767 have not the logico-mathematical grasp to defend it in any depth.  
3768 That has been done, in particular by Geoffrey Hellman. (T&T,  
3769 117)

3770 The reader who wants a thorough discussion of modal structuralism should  
3771 consult Hellman (1989). For now, a back-of-the-envelope summary will be  
3772 sufficient. The modal structuralist claims that the theory of the natural num-  
3773 bers is not a description of some particular sequence of entities; rather, the  
3774 theory concerns *all possible* models of the Peano axioms.<sup>5</sup>

3775 For example, when a mathematician asserts that there are infinitely many  
3776 prime numbers, what is really meant is something like this:<sup>6</sup>

3777 It is necessary that, in any model of the Peano axioms, the domain  
3778 contains infinitely many prime elements.

3779 Notice that this statement does not imply that there is a model of the Peano  
3780 axioms somewhere in spacetime.

3781 When a modal structuralist mathematician asserts that, necessarily, every  
3782 model of the Peano axioms contains infinitely many prime elements, she will  
3783 of course wish to rule out the suggestion that this is true “vacuously”—that  
3784 is, simply because models of the Peano axioms are *impossible*. Thus, a modal  
3785 structuralist mathematician will claim that models of the Peano axioms are  
3786 possible.

3787 Hellman discusses in some detail how to extend this approach beyond num-  
3788 ber theory, and into applied mathematics. We need not look into these details.  
3789 For our purposes, the key point is that the appeal to modal structuralism al-  
3790 lows Armstrong to say that “ $10^{10^{10}}$  is even” is true (when properly interpreted)

5 The Peano axioms are the standard axioms in the theory of the natural numbers. Among them are such claims as “Zero is a number,” and “If  $x$  is any natural number,  $x + 0 = x$ .”

6 The version of modal structuralism that I so quickly sketch here is hermeneutic rather than revolutionary (for this distinction, see Burgess, J.P. and Rosen 1997, 6–7). This is, I think, the correct interpretation of Armstrong’s position. For Hellman’s position, see (1998).

3791 without committing himself to the questionable thesis that  $10^{10^{10}}$  physical  
3792 objects exist.<sup>7</sup>

3793 The attractions of the approach are obvious, but Armstrong’s possibilism  
3794 brings with it a new problem. Armstrong was a truthmaker maximalist—he  
3795 believed that every true proposition has a “truthmaker,” that is, an entity in  
3796 the spacetime world which is sufficient (and perhaps more than sufficient)  
3797 to explain the truth of the proposition.<sup>8</sup> Armstrong was thus stuck with the  
3798 formidable task of identifying truthmakers for complex modal truths like  
3799 those described above. It is my contention that Armstrong did not succeed at  
3800 this task, as I shall explain in the next three sections.<sup>9</sup>

### 3803 Armstrong’s Entailment Principle

3802 Before we look at Armstrong’s discussion of truthmaking and modality, we  
3803 must consider his ENTAILMENT PRINCIPLE, which is crucial to his account.  
3804 Some notation will be helpful: I will put a sentence between angled brackets  
3805 to represent the corresponding proposition. For example,  $\langle \text{Sam is dancing} \rangle$  is  
3806 the proposition that Sam is dancing. Here is the ENTAILMENT PRINCIPLE, as  
3807 it is formulated in *Sketch for a Systematic Metaphysics*:

3808 ENTAILMENT PRINCIPLE (SSM VERSION). If  $\alpha$  entails  $\beta$ , then any  
3809 truthmaker for  $\langle \alpha \rangle$  must be a truthmaker for  $\langle \beta \rangle$  too. (SSM, 65–66)

3810 For example, since  $\varphi$  entails  $\neg\neg\varphi$ , it is a consequence of Armstrong’s ENTAIL-  
3811 MENT PRINCIPLE that any truthmaker for  $\langle \varphi \rangle$  must also be a truthmaker for

7 Hellman’s modal structuralism involves second-order quantification, and it is worth thinking about how such quantifiers should be interpreted within Armstrong’s metaphysical system. One approach is to say that the second-order variables range over properties (including “second-rate” properties—see footnote 2). Some restriction of the usual comprehension axiom will be needed to accommodate Armstrong’s contention that there are no uninstantiated properties. For a version of modal structuralism that does not require second order quantification, see Berry (2018).

8 The parenthetical “and perhaps more than sufficient” is there to indicate that Armstrong’s was an *inexact* conception of truthmaking, to use Kit Fine’s terminology—see (2017).

9 Fox (1987) endorses a *purely modal* conception of truthmaking. According to Fox, T is a truthmaker for  $p$  just in case it is necessary that if T exists then  $p$  is true. On this approach, it is easy to identify truthmakers for purely mathematical truths. Since the truths of pure mathematics are necessary, given Fox’s purely modal conception of truthmaking, *anything whatever* is a truthmaker for any purely mathematical truth.

Armstrong himself vigorously rejected this approach, insisting that truthmakers must be *relevant* to the propositions they make true (T&T, 11). For more on this theme, see Cameron (2018).

3812  $\langle \neg\neg\varphi \rangle$ . In this example, the ENTAILMENT PRINCIPLE is plausible. However,  
 3813 there is an important objection to this formulation of the principle. As Restall  
 3814 (1996) has pointed out, this simple version of the ENTAILMENT PRINCIPLE  
 3815 conflicts with a popular and appealingly simple (though not undisputed)  
 3816 account of truthmaking and disjunction:

3817 DISJUNCTION PRINCIPLE. T makes true the proposition  $\langle \varphi \vee \psi \rangle$  if  
 3818 and only if T makes true  $\langle \varphi \rangle$ , or T makes true  $\langle \psi \rangle$ , or both.

3819 To see the conflict, consider the following argument:

3820 Let  $\varphi$  and  $\psi$  be any two true sentences, and suppose that T is a  
 3821 truthmaker for  $\langle \varphi \rangle$ . Then since  $\varphi$  entails  $(\psi \vee \neg\psi)$ , T must also be a  
 3822 truthmaker for  $\langle \psi \vee \neg\psi \rangle$ . By the DISJUNCTION PRINCIPLE, T must  
 3823 be a truthmaker either for  $\langle \psi \rangle$  or for  $\langle \neg\psi \rangle$ . But by hypothesis,  $\langle \psi \rangle$  is  
 3824 true so  $\langle \neg\psi \rangle$  is false and so  $\langle \neg\psi \rangle$  has no truthmakers. So T must be  
 3825 a truthmaker for  $\langle \psi \rangle$ .

3826 This little argument appears to show for any two true sentences  $\varphi$  and  $\psi$ , any  
 3827 truthmaker for  $\langle \varphi \rangle$  is also a truthmaker for  $\langle \psi \rangle$ —a result which completely  
 3828 trivializes truthmaker theory.

3829 In *Truth and Truthmakers*, Armstrong gives a more sophisticated version  
 3830 of the ENTAILMENT PRINCIPLE which is not subject to the same objection:

3831 ENTAILMENT PRINCIPLE (T&T VERSION). If  $\alpha$  entails\*  $\beta$ , then any  
 3832 truthmaker for  $\langle \alpha \rangle$  must be a truthmaker for  $\langle \beta \rangle$  too. (T&T, 10)

3833 Here, entailment\* is some non-classical entailment relation, to be specified.  
 3834 By insisting that  $\varphi$  need not entail\*  $(\psi \vee \neg\psi)$ , we can maintain a version of  
 3835 the ENTAILMENT PRINCIPLE without having to conclude, absurdly, that all  
 3836 propositions expressed by true sentences have the same truthmakers.<sup>10</sup>

3837 We have seen that the SSM VERSION of the ENTAILMENT PRINCIPLE con-  
 3838 flicts with the DISJUNCTION PRINCIPLE, and that one can maintain the DIS-  
 3839 JUNCTION PRINCIPLE by endorsing the T&T VERSION instead. If a truthmaker  
 3840 theorist wishes instead to maintain the simpler, SSM VERSION of the ENTAIL-

10 It is not easy to say exactly what entailment\* is. Restall (1996) has proposed that entailment\* “is nearly, but not quite, the first degree entailment of relevant logic.” Linnebo (2022) has suggested that entailment\* includes first-order intuitionistic entailment without identity. Thankfully, we need not settle this question here.

3841 MENT PRINCIPLE, she may choose to reject the **DISJUNCTION PRINCIPLE**—and  
 3842 indeed philosophers have presented *independent* reasons for rejecting this  
 3843 principle.<sup>11</sup> We need not settle this dispute here. Suffice it to say that appeals  
 3844 to the **SSM VERSION** of the **ENTAILMENT PRINCIPLE** are subject to dispute.<sup>12</sup>

#### 3844 **4 Armstrong on Truthmaking and Possibility**

3846 Having briefly looked at the **ENTAILMENT PRINCIPLE** we are ready to consider  
 3847 Armstrong’s account of truthmaking and modality. Let’s start with possibility.  
 3848 Suppose that the sentence  $\varphi$  expresses a contingently true proposition; what  
 3849 then are the truthmakers for  $\langle \diamond\varphi \rangle$  and  $\langle \diamond\neg\varphi \rangle$ ?

3850  $\langle \diamond\varphi \rangle$  is comparatively straightforward. Since  $\langle \varphi \rangle$  is true, Armstrong argued,  
 3851 it must have a truthmaker, T. Since  $\varphi$  entails  $\diamond\varphi$ , T will be a truthmaker for  
 3852  $\langle \diamond\varphi \rangle$  as well, by the **ENTAILMENT PRINCIPLE**.<sup>13</sup>

3853  $\langle \diamond\neg\varphi \rangle$  is rather more difficult. Armstrong introduced his “possibility prin-  
 3854 ciple” (T&T, 84) to deal with the problem:

3855 **POSSIBILITY PRINCIPLE.** If  $\langle \varphi \rangle$  is a contingent truth and T is a  
 3856 truthmaker for  $\langle \varphi \rangle$ , then T is a truthmaker for  $\langle \diamond\neg\varphi \rangle$ .

3857 The principle is not attractive on its face. As Pawl (2010) has pointed out,  
 3858 Armstrong’s being legged is a truthmaker for  $\langle \text{Someone has legs} \rangle$ , but it is  
 3859 hardly plausible that Armstrong’s being legged is a truthmaker for  $\langle \text{Possibly,}$   
 3860  $\text{nobody has legs} \rangle$ . But Armstrong claimed that the **POSSIBILITY PRINCIPLE** is  
 3861 a consequence of the **ENTAILMENT PRINCIPLE**. He presented the following  
 3862 argument (T&T, 84, notation slightly modified):

- 
- |   |              |
|---|--------------|
| (1) T is a truthmaker for $\langle \varphi \rangle$ . | (Assumption) |
| (2) $\langle \varphi \rangle$ is contingent.          | (Assumption) |
- 

11 For discussion of the **DISJUNCTION PRINCIPLE**, see Rodriguez-Pereyra (2006) and López de Sa (2009).

12 Why did Armstrong give these two different versions of the **ENTAILMENT PRINCIPLE**, in books written in the same decade? My hypothesis is that T&T contains Armstrong’s preferred formulation of the **ENTAILMENT PRINCIPLE**, and that the version in SSM (a much shorter, easier work) is a simplification.

13 It is not entirely straightforward that this application of the **ENTAILMENT PRINCIPLE** is correct. While it is clear that  $\varphi$  entails  $\diamond\varphi$ , it is perhaps not so clear that  $\varphi$  entails\*  $\diamond\varphi$ . We can let this point slide, however, because there are much more serious objections to Armstrong’s position, as we shall see.

- 
- |     |   |   |
|-----|---|---|
| (3) | $\langle \varphi \rangle$ entails $\langle \Diamond \neg \varphi \rangle$ . | (From (2), and the nature of the contingency of propositions) |
| (4) | T is a truthmaker for $\langle \Diamond \neg \varphi \rangle$ .             | (From (1), (3) and the ENTAILMENT PRINCIPLE)                  |
- 

3863 As Armstrong later recognized, this argument is fallacious. The error is in  
 3864 step (3): in no standard modal logic is it true that  $\langle \varphi \rangle$  entails  $\langle \Diamond \neg \varphi \rangle$  (special  
 3865 cases aside).

3866 So the **POSSIBILITY PRINCIPLE** is implausible on its face, and the argument  
 3867 Armstrong gave for it is unconvincing. I think we should conclude that the  
 3868 principle should be rejected.<sup>14</sup>

3869 Happily, Armstrong also offered another and more attractive account of  
 3870 truthmaking and possibility. The idea is this. Suppose that we have (separately)  
 3871 two slices of bread, fifteen slices of cheese, and two slices of tomato. These  
 3872 things could have constituted a cheese and tomato sandwich—although in  
 3873 fact they don't. Plausibly, they together form a truthmaker for the proposition  
 3874 that a cheese and tomato sandwich *could* exist. Armstrong wrote:

3875 Consider, in particular, the cases where the entities in question  
 3876 do not exist, where they are *mere* possibilities. It is, let us sup-  
 3877 pose, true that  $\langle$ it is possible that a unicorn exists $\rangle$ . What then  
 3878 is a minimal truthmaker for this truth? The obvious solution is  
 3879 combinatorial. The non-existent entity is some non-existent (but  
 3880 possible) combination out of elements that do exist. The phrase  
 3881 “non-existent combination” may raise eyebrows. Am I commit-  
 3882 ting myself to a Meinongian view? No, I say. The *elements* of the  
 3883 combination are, I assert, the only truthmakers that are needed  
 3884 for the truth that this combination is possible. (T&T, 91–92)

3885 I think that this is a very attractive account of what the truthmakers are for  
 3886 *some* truths about possibility (including truths about unicorns and tomato  
 3887 sandwiches).<sup>15</sup> However, it is doubtful that the combinatorial approach pro-  
 3888 vides us with sufficient truthmakers for all the possibility claims made by  
 3889 the modal structuralist mathematician. The modal structuralist will assert

---

14 Armstrong (2007) recognized the error in his argument for the **POSSIBILITY PRINCIPLE** and went on to offer a *new* argument for it. For criticism of this later argument, see Pawl (2010).

15 For some criticisms of Armstrong's “combinatorialist” theory of modality, see Wang (2013).



3890 that second-order ZFC could have had a model, but it seems unlikely that  
 3891 such a model could be created by recombining physical objects, because it  
 3892 seems unlikely that there are *enough* physical objects to go around. If, for  
 3893 example, there are only  $\beth_3$  physical objects, we will not by combining them  
 3894 be able to produce a set with  $\beth_4$  elements, but *every* model of second-order  
 3895 ZFC contains sets with  $\beth_4$  elements—and indeed much larger sets to boot.  
 3896 The problem of size has reemerged in a new form.

## 3897 **5 Armstrong on Truthmaking and Necessity (Part 1: Truth** 3898 **and Truthmakers)**

3899 Let's turn to Armstrong's discussion of propositions about necessity. Since our  
 3900 concern is Armstrong's philosophy of mathematics, we need not discuss all  
 3901 of Armstrong's views about truthmaking and necessity. Instead, we'll focus  
 3902 on what he had to say about truthmakers for the theorems of mathematics.  
 3903 Armstrong's discussions of this topic in *Truth and Truthmakers* and *Sketch*  
 3904 *for a Systematic Metaphysics* are very different. In this section, I'll consider  
 3905 chapter eight of *Truth and Truthmakers*, leaving the later book until section 7.

3906 Armstrong (T&T, 99, 111) suggested that the numbers themselves constitute  
 3907 truthmakers for some arithmetical truths. For example, seven, five and twelve  
 3908 may together form a truthmaker for  $\langle 7 + 5 = 12 \rangle$ .<sup>16</sup> For Armstrong, so long as  
 3909 seven, five and twelve exist they *must* be related in this way, and so nothing  
 3910 *beyond* their existence is needed to explain their being so related. This relation  
 3911 between the three numbers is "internal" to them.

3912 This is an important idea, and I will return to it in section 6. But this is not  
 3913 on its own a complete solution to the problem at hand. Consider for example  
 3914 the proposition  $\langle \beth_\omega + \beth_\omega = \beth_\omega \rangle$ . This is a theorem of orthodox mathematics,  
 3915 and so Armstrong would surely accept that the proposition is true, when given  
 3916 its proper modal interpretation. But what is its truthmaker? Surely  $\beth_\omega$  itself  
 3917 can be a truthmaker only if it exists. However, as we saw in section 1 and  
 3918 section 2, it is doubtful for Armstrong that  $\beth_\omega$  exists.<sup>17</sup>

3919 Later in the chapter, Armstrong discussed analytic truths. He wrote:

16 A variant on this suggestion (T&T, 98) is that any truthmaker for  $\langle 7 \text{ exists} \rangle$ ,  $\langle 5 \text{ exists} \rangle$  and  $\langle 12 \text{ exists} \rangle$  will also be a truthmaker for  $\langle 7 + 5 = 12 \rangle$ .

17 Of course, for Armstrong, " $\beth_\omega$  exists" is true when given a suitable modal reinterpretation—but this is only because, so interpreted, the sentence doesn't actually assert the existence of  $\beth_\omega$ !

3920 A traditional view, which has many supporters, is that [analytic]  
 3921 truths are true solely in virtue of the meanings of the terms in  
 3922 which they are expressed. (T&T, 109)

3923 Armstrong went on to say that “[t]he phrase ‘in virtue of’ inevitably suggests  
 3924 truthmakers.” So Armstrong proposed that if a sentence *S* is analytic, the  
 3925 proposition it expresses is made true by the meanings of the words in *S*. For  
 3926 example, ⟨A father is a male parent⟩ is made true by the meanings of “a,”  
 3927 “father,” “is,” “a,” “male” and “parent.”

3928 Now Armstrong suggested—somewhat tentatively—that statements in  
 3929 mathematics about what is necessary are analytic.<sup>18</sup> On this view, the mean-  
 3930 ings of mathematical terms make true all such statements.

3931 I do not dismiss completely the claim that mathematical truths are ana-  
 3932 lytic.<sup>19</sup> However, Armstrong’s version of this thesis is insufficient to solve the  
 3933 problem at hand. Let *p* be some true proposition from pure mathematics. Arm-  
 3934 strong believed that the theorems of pure mathematics are *necessarily* true.  
 3935 So we can ask what the truthmaker for *p* would have been, had there been no  
 3936 language-users. How might Armstrong reply? Surely it is not adequate to say  
 3937 that the meanings of English words would have been the truthmakers—for  
 3938 English words would not have existed in the absence of English speakers.<sup>20</sup>  
 3939 If Armstrong replies that *p* would not have had a truthmaker, he would be  
 3940 stuck with the surely unwanted conclusion that it is possible for a proposition  
 3941 to be true without a truthmaker. And so the proper Armstrongian conclusion  
 3942 is that *p* would have had a different set of truthmakers, had there been no  
 3943 language-users. But then we are left with the question of what these truth-  
 3944 makers would have been—and until this question is answered, Armstrong’s  
 3945 account is incomplete.

---

18 Armstrong wrote: “There may be something mechanical, something purely conceptual, purely semantic, in the deductive following-out of proofs of the existence of the possible. (See the account of analytic truth to come in 8.9.)” (T&T, 102). Note that this quotation is from the chapter on *necessary* truths in T&T. So I take it that what Armstrong is (tentatively) suggesting here is that truths in mathematics about what is necessary are analytic.

19 When Armstrong says that “a traditional view” is that analytic truths are “true solely in virtue of the meanings of the terms in which they are expressed,” his wording seems to derive from the introduction to Ayer (1946). I think that few philosophers of mathematics today would defend Ayer’s view in all its details. However, there are still philosophers who endorse views which resemble Ayer’s position in important respects. See for example Rayo (2013).

20 Perhaps some philosophers will insist that words (or their meanings) are necessarily existing abstract objects. However, I take it that Armstrong would not take this line. As we’ve seen, Armstrong rejected necessary abstracta.

3944 **6 Armstrong on Truthmaking and Necessity (Part 2: Sketch**  
 3947 **for a Systematic Metaphysics)**

3948 By the time he wrote *Sketch for a Systematic Metaphysics*, Armstrong had  
 3949 decided to reject his earlier suggestion that mathematical truths are analytic,  
 3950 saying that such a view implies that mathematics is “too arbitrary or conven-  
 3951 tional” (SSM, 91). But he suggested an alternative approach, which we will  
 3952 now consider.

3953 We should begin by looking at Armstrong’s metaphysics of law. Armstrong  
 3954 claimed that a law is a relation between properties (SSM, 35). Here is a toy  
 3955 example. Suppose that it is a law that being dehydrated causes headaches. For  
 3956 Armstrong, this means that a certain relation (viz.  $\mathcal{N}$ , the nomic relation) ob-  
 3957 tains between two properties (viz. the property *is dehydrated*, and the property  
 3958 *has a headache*). Armstrong would symbolize this as follows:<sup>21</sup>

3959  $\mathcal{N}(\textit{is dehydrated}, \textit{has a headache})$

3960 Now Armstrong claimed that laws have “instantiations.” Our law, for example,  
 3961 is instantiated whenever someone is dehydrated and, consequently, has a  
 3962 headache. A law, on this view, is itself a property. And we have already seen  
 3963 that Armstrong was happy to posit “immanent” properties. This led Armstrong  
 3964 to the view that every law is instantiated. He wrote:

3965 If laws are a species of universal, then, according to me at least,  
 3966 they have to be instantiated at some place and time. Well, we talk  
 3967 of laws being instantiated, do we not? (The points where the laws  
 3968 are ‘operative’.) So this instantiation of laws is the instantiation  
 3969 of a special sort of universal. (Note that this would require ev-  
 3970 ery law to be somewhere instantiated in space-time.) [...] One  
 3971 consequence of this is that there cannot be laws that are never  
 3972 instantiated. (SSM, 41)

3973 Now Armstrong suggested that it is certain mathematical and logical laws  
 3974 which make true the necessities of mathematics.

3975 He began his discussion of this proposal by appealing to his ENTAILMENT  
 3976 PRINCIPLE, arguing that truthmakers for the axioms of a mathematical theory

---

21 I allow myself here to omit some of the finer details of Armstrong’s account—in particular, I do not mention “state of affairs types” (SSM, 28–40).

3977 must also be truthmakers for the theorems (SSM, 90). This manoeuvre is  
 3978 suspect. While it is known that the theorems of orthodox mathematics are  
 3979 entailed by the axioms (that's what makes them theorems, after all) it is far  
 3980 from clear that the axioms entail\* the theorems.<sup>22</sup>

3981 And it gets worse. To complete his account, Armstrong still needed to  
 3982 specify truthmakers for the axioms of our mathematical theories. To do this,  
 3983 he appealed to his theory of laws:

3984 We do, of course, have to recognize that introducing the Entail-  
 3985 ment Principle drives us back to consider the axioms from which  
 3986 mathematical systems are developed. [...] I suggest that we should  
 3987 postulate laws in logic and mathematics (non-contradiction, ex-  
 3988 cluded middle in logic, Peano's axioms for number, or whatever  
 3989 laws logicians and mathematicians wish to postulate). In the light  
 3990 of the nature of proof just argued for we might suggest that such  
 3991 laws might be all we needed to postulate in the way of an ontology  
 3992 for logical and mathematical entities. (SSM, 90–91)

3993 It is not credible, as Armstrong suggests here, that the Peano axioms are “laws”  
 3994 in Armstrong's sense. For example, one of the Peano axioms states that the  
 3995 natural numbers are *unending* in the sense that every natural number has  
 3996 a successor.<sup>23</sup> For Armstrong, this statement may not be true when taken at  
 3997 face value. For Armstrong, as we've seen, the existence of very large natural  
 3998 numbers is doubtful, and it is at least possible that there is a *largest* natural  
 3999 number, which has no successor. To circumvent this point, Armstrong will  
 4000 presumably insist on a modal reinterpretation of the axiom. On this view, the  
 4001 axiom, properly interpreted, states that, necessarily, every model of the Peano  
 4002 axioms is unending.

4003 The corresponding Armstrongian law would then have to be:

- 
- 22 Suppose, for example, that Linnebo (2022) is correct and entailment\* coincides with intuitionistic entailment. Then consider some statement  $\varphi$  which is provable classically but not intuitionistically from the prevailing axioms. (For example,  $\varphi$  might be  $(\psi \vee \neg\psi)$ , where  $\psi$  is some statement independent of the prevailing axioms.) Assuming, with Armstrong, that the inferences of classical logic are all truth-preserving, and the axioms of orthodox mathematics are true, we can conclude that  $\varphi$  is true. However, because it is not entailed\* by the prevailing axioms, we cannot identify a truthmaker for it using the proposed method.
- 23 The “successor” of a natural number is the number that comes *immediately after it*, when the natural numbers are arranged in the customary fashion. So for example the successor of nineteen is twenty.

4004 *N(is a model of Peano arithmetic, is unending)*

4005 But this is problematic, because Armstrong believed that all laws are instanti-  
4006 ated in the physical world—and it is far from clear that this law is instantiated.  
4007 It may be that the physical world is finite. However, every model of Peano  
4008 arithmetic is infinite. And so it may be that there are no physical models of  
4009 Peano arithmetic, in which case the above-mentioned law is uninstantiated.

4010 We might be able to avoid this problem by arguing on empirical grounds  
4011 that there are infinitely many physical objects. For example, we might appeal  
4012 to the common (though admittedly contested) assumption in physics that  
4013 space is infinitely divisible. However, the problem that I have just described  
4014 will reassert itself when we turn our attention from Peano arithmetic to other  
4015 branches of mathematics which posit a greater number of entities. The most  
4016 extreme case is set theory. Any model of second-order ZFC would have to have  
4017 a truly vast domain, containing  $\aleph_\omega$  elements and more. There is no empirical  
4018 reason to think that there exist that many physical objects. So we are left with  
4019 the conclusion that the Armstrongian laws corresponding to the axioms of  
4020 set theory are uninstantiated.

4021 To avoid these problems, an Armstrongian would have to list a number of  
4022 basic principles for mathematics which express laws that *are* instantiated in  
4023 the physical world, and argue that they entail\* the truths of mathematics. I  
4024 don't know that this is impossible, but it is far from obvious that it can be  
4025 done. And even if it *could* be done, the problem of identifying truthmakers  
4026 for facts about what is possible would remain.

## 4027 **An Alternative Approach**

4028 Let's review. Armstrong believed that mathematical entities are located within  
4029 the physical world. For example, wherever there is a pair of things, there is the  
4030 number two. However, Armstrong realized that the physical world is not large  
4031 enough to accommodate all the entities posited by modern pure mathematics.  
4032 So he adopted a modal interpretation of mathematics. For Armstrong, pure  
4033 mathematics tells us not about what is, but about what could be and what  
4034 must be. However, Armstrong believed that every truth has a truthmaker  
4035 within the physical world, and so he was left with the unenviable task of  
4036 identifying truthmakers for modal truths within the physical world. I have

4037 argued that he did not succeed.<sup>24</sup> In this final section, I would like to put  
 4038 forward an alternative approach—an approach which will, I hope, appeal to  
 4039 those impressed by Armstrong’s metaphysical system.<sup>25</sup>

4040 Armstrong accepted a version of the methodological principle known as  
 4041 Occam’s razor. He rejected mathematical Platonism largely for this reason.  
 4042 A “Platonic realm of numbers,” he wrote, is an “ontological extravagance”  
 4043 (T&T, 100).<sup>26</sup> However, Armstrong did not use his razor to excise supervenient  
 4044 entities. Supervenient entities, he thought, are an “ontological free lunch.”  
 4045 For example, he did not think that universalism in mereology is objectionably  
 4046 unparsimonious:

4047       Whatever supervenes or, as we can also say, is entailed or necessi-  
 4048       tated, is not something ontologically additional to the subvenient,  
 4049       or necessitating, entity or entities. [...] The terminology of “noth-  
 4050       ing over and above” seems appropriate to the supervenient. [...] If  
 4051       the supervenient is not something ontologically additional, then  
 4052       this gives charter to, by exacting a low price for, an almost entirely  
 4053       permissive mereology. Do the number 42 and the Murrumbidgee  
 4054       River form a mereological whole? [...] The whole, if it exists, is  
 4055       certainly a strange and also an uninteresting object. But if it super-  
 4056       venes on its parts, and if as a consequence of supervening it is not  
 4057       something more than its parts, then there seems no objection to  
 4058       recognizing the whole. So in this essay permissive mereology, un-  
 4059       restricted mereological composition, is embraced. (1997, 12–13)

4060 On an uncharitable interpretation of this passage, Armstrong’s view was that  
 4061 if the existence of *x* necessitates the existence of *y*, then *y* is “nothing over  
 4062 and above” *x*. But this is hardly plausible. Perhaps God exists necessarily, but  
 4063 it would be grossly immodest for me to claim that God is nothing over and  
 4064 above me. Perhaps I could not have had different parents, in which case my  
 4065 existence necessitates theirs. But they would quite properly take exception to  
 4066 the suggestion that they are nothing ontologically additional to me.<sup>27</sup>

24 It is worth noting in passing that Armstrong’s theory of *propositions* was problematic in rather similar ways. On this point, see McDaniel (2005).

25 For a very different approach, see Read (2010).

26 Armstrong also had epistemological reasons for rejecting Platonism (SSM, 2). For lack of space, I do not discuss epistemology in this paper.

27 For a more detailed discussion of these points, see Schulte (2014).

4067 Cameron (2008) has suggested a more promising way of developing Arm-  
 4068 strong's idea that supervenient entities are "free." To put it briefly, Cameron's  
 4069 proposal is as follows. Compare the following two propositions:

4070  $\langle m \text{ exists.} \rangle$  (where  $m$  is a marriage, between Ashni and Ben)  
 4071  $\langle e \text{ exists.} \rangle$  (where  $e$  is an electron)

4072 The former proposition is made true by certain patterns of human activity—  
 4073 involving perhaps Ashni, Ben, a registrar, some pieces of paper, and some  
 4074 metal rings. Ashni and Ben's marriage is a *derivative* entity: its existence is  
 4075 explained by facts about things other than itself. The electron  $e$  is not derivative.  
 4076 The electron's existence is not explained by facts about other things;  $e$  itself is  
 4077 the only truthmaker for the proposition  $\langle e \text{ exists.} \rangle$ .<sup>28</sup>

4078 More generally, Cameron's proposal is this. When  $x$  is fundamental, the  
 4079 only truthmaker for  $\langle x \text{ exists} \rangle$  is  $x$  itself. When  $x$  is derivative,  $\langle x \text{ exists} \rangle$  has a  
 4080 truthmaker other than  $x$  itself.<sup>29</sup>

4081 Cameron adds that it is derivative entities *in this sense* that are an "ontolog-  
 4082 ical free lunch," to use Armstrong's phrase. In effect, Cameron replaces the  
 4083 familiar slogan "Do not multiply entities beyond necessity" with a variant:  
 4084 "Do not multiply *fundamental* entities beyond necessity". Since mereological  
 4085 compounds are non-fundamental, Cameron infers, mereological universalism  
 4086 is not objectionable on grounds of parsimony.<sup>30</sup>

4087 Cameron briefly suggests an application of this idea to impure set theory.  
 4088 He proposes that an impure set is "nothing over and above" its elements, so  
 4089 there is no objection on grounds of parsimony to positing all those impure  
 4090 sets that can be built up from *basic elements* whose existence can already  
 4091 be established. On this view, there is no need to re-interpret set theory in a  
 4092 "possibilist" manner. We maintain that all the sets posited by set theorists really

28 Sharp-eyed readers will note that in this section I assume an *explanatory* conception of truthmaking, according to which, when  $T$  is a truthmaker for  $p$ ,  $T$  *explains the truth of*  $p$ . For discussion, see Cameron (2018).

29 I have actually modified Cameron's proposal in a small way. Cameron's view is that when  $x$  is derivative,  $x$  is *not* a truthmaker for  $\langle x \text{ exists} \rangle$ . I find this claim puzzling (How could  $x$  fail to make true  $\langle x \text{ exists} \rangle$ ?) and since it is inessential to my argument, I omit it.

30 Suppose that  $a$  and  $b$  are fundamental objects, and that  $(a + b)$  is their mereological sum. According to Cameron  $a$  and  $b$  *collectively* make true  $\langle (a + b) \text{ exists} \rangle$ . For Cameron, this proposition has no *single* truthmaker; rather there are some things which *together* make the proposition true. This is a subtlety of Cameron's view which I ignore in the main text, for simplicity.

4093 do exist, although they are not fundamental. Let's develop this Cameronian  
4094 proposal in more detail.

4095 Why does the set {Jill, Joe} exist? I suggest that it exists because Jill exists,  
4096 and because Joe exists—and that is all. Nothing more is needed. And so, I sug-  
4097 gest, any truthmaker for  $\langle$ Jill exists $\rangle$  and  $\langle$ Joe exists $\rangle$  will also be a truthmaker  
4098 for  $\langle$ A exists $\rangle$ , where A is {Jill, Joe}. More generally:

4099 (1) If T is a truthmaker for  $\langle$ x exists $\rangle$ , for each x in a non-empty set X, then  
4100 T is a truthmaker for  $\langle$ X exists $\rangle$  also.<sup>31</sup>

4101 So much for propositions about the existence of sets. But a complete truth-  
4102 maker theoretic account of the sets will also include an account of what the  
4103 truthmakers are for other propositions, including propositions about the iden-  
4104 tity and distinctness of sets, and propositions about what is an element of  
4105 what.

4106 Let's start with identity. Suppose that someone asks us why Joe is identical  
4107 to Joe—that is, we are asked why Joe is identical with himself. This is a very  
4108 peculiar question. The best answer to it that I can come up with goes like  
4109 this. For Joe to bear the identity relation to himself, it suffices that he exists.  
4110 Self-identity is not some additional characteristic that requires further  
4111 explanation. Joe exists, and so he is self-identical. And that is that. If this is  
4112 right, I suggest, any truthmaker for the proposition  $\langle$ Joe exists $\rangle$  must also be a  
4113 truthmaker for  $\langle$ Joe = Joe $\rangle$ . In general, any truthmaker for  $\langle$ x exists $\rangle$  will also  
4114 be a truthmaker for  $\langle$ x = x $\rangle$ .<sup>32</sup>

4115 Something similar is plausible in the case of non-identity. If, bizarrely, we  
4116 are asked why it is that Jill is not identical with Joe—if we are asked why they  
4117 are two people and not one—all we can say in reply is that to be non-identical

---

31 What about the empty set? One might be tempted to avoid the problem by denying that the empty set exists. This proposal, as Hazen (1991) argues, is less radical than it might first seem, and Armstrong did in some places express scepticism about the empty set (T&T, 114). However, given Armstrong's usual hostility towards philosophically motivated reforms to standard mathematical practice, I think it desirable, from an Armstrongian point of view, to preserve the empty set. So here is an alternative, inspired by Kit Fine's well-known discussion of zero-grounding (2012).

Armstrong generally supposed that a truthmaker will always be a *single thing*. But we might want to allow that a proposition can be made true by two things acting in concert, or three things, or four things, or more. For example, we might say that *a*, *b* and *c* collectively make true  $\langle$ {*a*, *b*, *c*} exists $\rangle$ . Taking this line of thought still further, we could argue that in some unusual cases a proposition is made true by *zero* things; as we might put it, such propositions are *trivially* made true. On this view, we may say that  $\langle$  $\emptyset$  exists $\rangle$  is trivially made true.

32 For more detailed discussion of the question of how truths of identity are to be explained, see Burgess, A. (2012) and Shumener (2017).



4118 the two Bidens need only exist. Jill exists. She is one person. Joe exists too.  
 4119 He is another. And that is all. There is nothing extra that Jill and Joe need to  
 4120 do or to be in order to be distinct—existing is enough. And so, I suggest, any  
 4121 truthmaker for  $\langle \text{Jill exists} \rangle$  and  $\langle \text{Joe exists} \rangle$  is a truthmaker also for  $\langle \text{Jill} \neq \text{Joe} \rangle$ .  
 4122 More generally, if  $x$  and  $y$  exist and are distinct, any truthmaker for  $\langle x \text{ exists} \rangle$   
 4123 and  $\langle y \text{ exists} \rangle$  must also be a truthmaker for  $\langle x \neq y \rangle$ .

4124 I want to recommend a similar treatment of the relations of membership  
 4125 and non-membership. If we are asked why Joe is an element of his singleton,  
 4126 there is nothing we can say except that, for this to be so, it suffices that Joe  
 4127 and his singleton exist. No more is needed. And if we are asked why Joe is  
 4128 not an element of  $\{\text{Jill}\}$ , we can say only that it is enough that Joe and  $\{\text{Jill}\}$   
 4129 exist. More generally, I suggest, if  $x$  is an element of  $Y$ , then any truthmaker  
 4130 for  $\langle x \text{ exists} \rangle$  and  $\langle Y \text{ exists} \rangle$  is a truthmaker too for  $\langle x \in Y \rangle$ . And if  $x$  is not  
 4131 an element of  $Y$ , though they both exist, any truthmaker for  $\langle x \text{ exists} \rangle$  and  
 4132  $\langle Y \text{ exists} \rangle$  is a truthmaker too for  $\langle x \notin Y \rangle$ .

4133 Let me put all of this in a rather different way. Let's say that a relation  $R$  is  
 4134 "strongly internal" if and only if the following condition is met: Necessarily,  
 4135 for any  $a$  and  $b$ , if  $a$  bears  $R$  to  $b$  then (1)  $a$  bears  $R$  to  $b$  at any world at which  $a$   
 4136 and  $b$  both exist, and (2) at every such world, any truthmaker for  $\langle a \text{ exists} \rangle$  and  
 4137  $\langle b \text{ exists} \rangle$  is also a truthmaker for  $\langle a \text{ bears } R \text{ to } b \rangle$ .<sup>33</sup> If  $R$  is a strongly internal  
 4138 relation and  $a$  bears  $R$  to  $b$ , then no explanation for this is required, beyond  
 4139 whatever is needed to account for the fact that the relata exist. My proposal is  
 4140 that the relations of identity, non-identity, membership and non-membership  
 4141 are strongly internal in this particular sense.<sup>34</sup> In summary:

- 4142 (1) If  $T$  is a truthmaker for  $\langle x \text{ exists} \rangle$ , for each  $x$  in a non-empty set  $X$ , then  
 4143  $T$  is a truthmaker for  $\langle X \text{ exists} \rangle$  also.

---

33 Armstrong said that a relation is internal if "given just the terms of the relation, the relation between them is necessitated" (T&T, 9). That is, given any relation  $R$ ,  $R$  is internal (in Armstrong's sense) just in case the following is necessary: For any  $a$  and  $b$ , if  $a$  bears  $R$  to  $b$ , then at every world at which  $a$  and  $b$  exist,  $a$  bears  $R$  to  $b$ .

Clearly, any strongly internal relation is also internal in Armstrong's sense.

The converse, however, is open to dispute. Suppose *arguendo* that God exists necessarily. Then the relation  $x$  and  $y$  are such that God exists is internal, in Armstrong's sense. But it is doubtful that this relation is strongly internal, for it is hardly plausible that any truthmaker for  $\langle \text{Joe exists} \rangle$  and  $\langle \text{Jill exists} \rangle$  must also be a truthmaker for  $\langle \text{God exists} \rangle$ .

34 Note that strongly internal relations need not be universals—they may be "second-rate" properties. In saying that non-membership is strongly internal, I do not assert that it is a genuine universal.

- 4144 (2) If a relation R is “strongly internal,” then whenever  $a$  bears R to  $b$ ,  
 4145 any truthmaker for  $\langle a \text{ exists} \rangle$  and  $\langle b \text{ exists} \rangle$  is also a truthmaker for  
 4146  $\langle a \text{ bears R to } b \rangle$ .  
 4147 (3) The relations of identity, non-identity, membership and non-  
 4148 membership are strongly internal.

4149 Let’s take this further. Suppose for example that T is a truthmaker for the  
 4150 proposition  $\langle \text{Jill exists} \rangle$ . Then by (1), T is also a truthmaker for each of these  
 4151 propositions:

- 4152  $\langle S_1 \text{ exists} \rangle$ , where  $S_1 = \{\text{Jill}\}$   
 4153  $\langle S_2 \text{ exists} \rangle$ , where  $S_2 = \{\{\text{Jill}\}\}$   
 4154  $\langle S_3 \text{ exists} \rangle$ , where  $S_3 = \{\{\{\text{Jill}\}\}\}$   
 4155  $\vdots$

4156 By (1) again, T is also a truthmaker for  $\langle S_\omega \text{ exists} \rangle$ , where  $S_\omega$  is the set  
 4157  $\{\text{Jill}, S_1, S_2, S_3, \dots\}$ .

4158 By (2) and (3), T is a truthmaker too for various propositions about the  
 4159 relations among these sets, propositions like  $\langle \text{Jill} \neq S_1 \rangle$ ,  $\langle S_1 = S_1 \rangle$ ,  $\langle S_1 \neq S_\omega \rangle$ ,  
 4160  $\langle S_\omega = S_\omega \rangle$ ,  $\langle \text{Jill} \in S_1 \rangle$ ,  $\langle S_1 \in S_\omega \rangle$ , and  $\langle S_1 \notin \text{Jill} \rangle$ .

4161 We can go further still, into the uncountable. Given our account, T will  
 4162 be a truthmaker for  $\langle S^* \text{ exists} \rangle$ , where  $S^*$  is the set of non-empty subsets of  
 4163  $S_\omega$ .  $S^*$  is an uncountable set. And T will be a truthmaker for  $\langle S^{**} \text{ exists} \rangle$ ,  
 4164 where  $S^{**}$  is the set of non-empty subsets of  $S^*$ —a set even larger than  $S^*$ .  
 4165 And proceeding in this way, we can locate in the physical world truthmakers  
 4166 for propositions concerning sets at all levels of the vertiginous set-theoretic  
 4167 hierarchy, including sets of arbitrarily high cardinality.

4168 And what of Armstrong’s claim that all entities exist “somewhere, some-  
 4169 when” (SSM, 15)? Well, some readers may find it edifying to insist that a set  
 4170 is located wherever its elements are.<sup>35</sup> On this view, you are co-located with  
 4171 your singleton, and its singleton, and *its* singleton, and so on *ad infinitum*. I  
 4172 offer no objection to this proposal. But I find it hard to *justify* the  
 4173 claim that sets have spatial locations, and more importantly it seems to me  
 4174 that we need not endorse this claim to earn the title “naturalist.” We insist


35 Maddy (1990) defends this view. I lack the space for a thorough treatment of Maddy’s approach, but I would like to note in passing that Maddy’s version involves *reforms* to standard set theory: to be specific, Maddy identifies individuals with their singletons (e.g., for Maddy, Socrates =  $\{\text{Socrates}\}$ ) and she eschews pure sets. The Armstrongian approach that I recommend preserves set theory in its usual form. (On the issue of pure sets, see footnote 31).

4175 that all fundamental objects are physical, and that all truths have physical  
4176 truthmakers—and this is naturalism enough.

4177 Back to the cardinal numbers. According to the current proposal, even  
4178 if the fundamental objects are rather few, nevertheless the sets are fantas-  
4179 tically numerous. This allows us to maintain Armstrong’s original account  
4180 of cardinal number without having to worry about the problem of size, and  
4181 without recourse to possibilism. Given the current proposal, for example,  $\aleph_\omega$   
4182 is instantiated in the hierarchy of sets, even if there are only finitely many  
4183 fundamental entities. If we add that it is not possible for there to be nothing,<sup>36</sup>  
4184 we are left with the conclusion that the cardinal numbers exist necessarily.

4185 Of course, a thorough truthmaker-theoretic account of mathematics would  
4186 also cover functions, complex numbers, matrices, ordinal numbers, graphs,  
4187 and all the other mathematical creatures. You will probably be relieved to  
4188 hear that I don’t intend to deal with all these topics now. It’s time for a cup of  
4189 tea, after all. But I hope that my discussion of sets and cardinals is sufficient  
4190 to motivate cautious optimism about Armstrongian naturalism—despite the  
4191 errors of detail that we have identified in Armstrong’s discussions of the  
4192 metaphysics of mathematics.\*

4193 Tom Donaldson

4194  0000-0001-6497-7038

4195 Simon Fraser University, Canada

4196 tmdonald@sfu.ca

## 4197 References

- 4198 ARMSTRONG, David M. 1989. *A Combinatorial Theory of Possibility*. Cambridge: Cam-  
4199 bridge University Press, doi:10.1017/cbo9781139172226.  
4200 —. 1997. *A World of States of Affairs*. Cambridge: Cambridge University Press, doi:10  
4201 .1017/cbo9780511583308.  
4202 —. 2004. *Truth and Truthmakers*. Cambridge Studies in Philosophy. Cambridge: Cam-  
4203 bridge University Press, doi:10.1017/cbo9780511487552.

36 Armstrong changed his mind on the question whether it is possible for there to be nothing. In Armstrong (1989, chap. 4, section IV) he claims that this is possible, but in T&T (105) he retracts the claim.

\* I would like to thank two reviewers at *Dialectica* for their hard work, and for their thorough and enormously helpful comments. I would also like to thank Jennifer Wang, Bahar Raiszadeh, Haktan Akcin, and Dans Waxman and Marshall for comments on drafts. (I’ve been working on this paper on and off for a long time; I apologize if I’m missing some acknowledgements!)

- 4204 —. 2007. “Truthmakers for Negative Truths and Truths of Mere Possibility.” in  
4205 *Metaphysics and Truthmakers*, edited by Jean-Maurice MONNOYER, pp. 99–104.  
4206 *Philosophische Analyse / Philosophical Analysis* n. 18. Heusenstamm  
4207 b. Frankfurt: Ontos Verlag, doi:10.1515/9783110326918.99.
- 4208 —. 2010. *Sketch for a Systematic Metaphysics*. Oxford: Oxford University Press, doi:10  
4209 .1093/acprof:oso/9780199590612.001.0001.
- 4210 AYER, Alfred Jules. 1946. *Language, Truth and Logic*. 2nd ed. London: Victor Gollancz.
- 4211 BARNES, Jonathan. 1985. “Aristotelian Arithmetic.” *Revue de philosophie ancienne* 3:  
4212 97–133. Reprinted, in revised form, in Barnes (2011, 334–363).
- 4213 —. 2011. *Method and Metaphysics*. Essays in Ancient Philosophy n. 1. Oxford: Oxford  
4214 University Press. Edited by Maddalena Bonelli.
- 4215 BERRY, Sharon E. 2018. “Modal Structuralism Simplified.” *Canadian Journal of Phi-*  
4216 *losophy* 48(2): 200–222, doi:10.1080/00455091.2017.1344502.
- 4217 BURGESS, Alexis. 2012. “A Puzzle about Identity.” *Thought* 1(2): 90–99, doi:10.1002/th  
4218 t3.14.
- 4219 BURGESS, John P. and ROSEN, Gideon. 1997. *A Subject with No Object: Strategies*  
4220 *for Nominalistic Interpretation of Mathematics*. Oxford: Oxford University Press,  
4221 doi:10.1093/0198250126.001.0001.
- 4222 CAMERON, Ross P. 2008. “Truthmakers and Ontological Commitment: or, How To  
4223 Deal with Complex Objects and Mathematical Ontology Without Getting into  
4224 Trouble.” *Philosophical Studies* 140(1): 1–18, doi:10.1007/s11098-008-9223-3.
- 4225 —. 2018. “Truthmakers.” in *The Oxford Handbook of Truth*, edited by Michael  
4226 GLANZBERG, pp. 333–354. Oxford Handbooks. New York: Oxford University  
4227 Press, doi:10.1093/oxfordhb/9780199557929.013.13.
- 4228 FINE, Kit. 2012. “Guide to Ground.” in *Metaphysical Grounding: Understanding*  
4229 *the Structure of Reality*, edited by Fabrice CORREIA and Benjamin  
4230 Sebastian SCHNIEDER, pp. 37–80. Cambridge: Cambridge University Press,  
4231 doi:10.1017/cbo9781139149136.002.
- 4232 —. 2017. “Truthmaker Semantics.” in *A Companion to the Philosophy of Language*,  
4233 edited by Bob HALE, Crispin WRIGHT, and Alexander MILLER, 2nd ed., pp. 556–577.  
4234 Blackwell Companions to Philosophy. Chichester: Wiley-Blackwell. First edition:  
4235 Hale and Wright (1997), doi:10.1002/9781118972090.ch22.
- 4236 FOX, John F. 1987. “Truthmaker.” *Australasian Journal of Philosophy* 65(2): 188–207,  
4237 doi:10.1080/00048408712342871.
- 4238 HALE, Bob and WRIGHT, Crispin, eds. 1997. *A Companion to the Philosophy of Lan-*  
4239 *guage*. Blackwell Companions to Philosophy. Oxford: Blackwell Publishers. Second  
4240 edition: Hale, Wright and Miller (2017).
- 4241 HALE, Bob, WRIGHT, Crispin and MILLER, Alexander, eds. 2017. *A Companion to the*  
4242 *Philosophy of Language*. 2nd ed. Blackwell Companions to Philosophy. Chichester:  
4243 Wiley-Blackwell. First edition: Hale and Wright (1997), doi:10.1002/978111897209  
4244 0.

- 4245 HAZEN, Allen Patterson. 1991. "Small Sets." *Philosophical Studies* 63(1): 119–123,  
4246 doi:10.1007/FBF00376001.
- 4247 HELLMAN, Geoffrey. 1989. *Mathematics without Numbers: Towards a Modal-Structural*  
4248 *Interpretation*. 11th ed. Oxford: Oxford University Press, doi:10.1093/0198240341  
4249 .001.0001.
- 4250 —. 1998. "Maoist Mathematics?" *Philosophia Mathematica* 6(3): 334–345, doi:10.1093/  
4251 philmat/6.3.334.
- 4252 KESSLER, Glenn. 1980. "Frege, Mill and the Foundations of Arithmetic." *The Journal*  
4253 *of Philosophy* 77(2): 65–79, doi:10.2307/2025431.
- 4254 LINNEBO, Øystein. 2022. "Generality Explained." *The Journal of Philosophy* 119(7):  
4255 349–379, doi:10.5840/jphil2022119725.
- 4256 LÓPEZ DE SA, Dan. 2009. "Disjunctions, Conjunctions, and their Truthmakers." *Mind*  
4257 118(470): 417–425, doi:10.1093/mind/fzp063.
- 4258 MADDY, Penelope. 1990. "Physicalist Platonism." in *Physicalism in Mathematics*, edited  
4259 by Andrew David IRVINE, pp. 259–289. The University of Western Ontario Series  
4260 in Philosophy of Science n. 45. Dordrecht: Kluwer Academic Publishers, doi:10.1  
4261 007/978-94-009-1902-0\_10.
- 4262 MCDANIEL, Kris. 2005. "Review of Armstrong (2004)." *Notre Dame Philosophical*  
4263 *Reviews*, 2005, <https://ndpr.nd.edu/reviews/truth-and-truthmakers/>.
- 4264 PAWL, Timothy. 2010. "The Possibility Principle and the Truthmakers for Modal  
4265 Truths." *Australasian Journal of Philosophy* 88(3): 417–428, doi:10.1080/0004  
4266 8400903193353.
- 4267 RAYO, Agustín. 2013. *The Construction of Logical Space*. Oxford: Oxford University  
4268 Press, doi:10.1093/acprof:oso/9780199662623.001.0001.
- 4269 READ, Stephen. 2010. "Necessary Truth and Proof." *Kriterion: Revista de Filosofia*  
4270 51(121): 47–67, doi:10.1590/S0100-512X2010000100003.
- 4271 RESTALL, Greg. 1996. "Truthmakers, Entailment and Necessity." *Australasian Journal*  
4272 *of Philosophy* 74(2): 331–340, doi:10.1080/00048409612347331.
- 4273 RODRIGUEZ-PEREYRA, Gonzalo. 2006. "Truthmaking, Entailment, and the Conjunction  
4274 Thesis." *Mind* 115(460): 957–982, doi:10.1093/mind/fzl957.
- 4275 ROSEN, Gideon. 1995. "Armstrong on Classes as States of Affairs." *Australasian Journal*  
4276 *of Philosophy* 73(4): 613–625, doi:10.1080/00048409512346971.
- 4277 SCHULTE, Peter. 2014. "Can Truthmaker Theorists Claim Ontological Free Lunches?"  
4278 *European Journal of Philosophy* 22(2): 249–268, doi:10.1111/j.1468-0378.2011.00  
4279 491.x.
- 4280 SHUMENER, Erica. 2017. "The Metaphysics of Identity: Is Identity Fundamental?"  
4281 *Philosophy Compass* 12(1): e12397, doi:10.1111/phc3.12397.
- 4282 SIMONS, Peter M. 1982. "Against the Aggregate Theory of Number." *The Journal of*  
4283 *Philosophy* 79(3): 163–167, doi:10.2307/2026072.
- 4284 WANG, Jennifer. 2013. "From Combinatorialism to Primitivism." *Australasian Journal*  
4285 *of Philosophy* 91(3): 535–554, doi:10.1080/00048402.2012.722114.

- 4286 YARNELLE, John Edward. 1964. *An Introduction to Transfinite Mathematics*. Boston,  
4287 Massachusetts: DC Heath and Company.

PROOF

4288

Published by *Philosophie.ch*

4289

Verein philosophie.ch

4290

Fabrikgässli 1

4291

2502 Biel/Bienne

4292

Switzerland

4293

[dialectica@philosophie.ch](mailto:dialectica@philosophie.ch)

4294

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

4295

ISSN 1746-8361

4296

ISBN 1234-5678

4297

This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

4298

4299

*Dialectica* is supported by the [Swiss Academy of Humanities and Social Sciences](#).

4300

4301

#### Abstracting and Indexing Services

4302

The journal is indexed by the Arts and Humanities Citation Index, Current Contents, Current Mathematical Publications, Dietrich's Index

4303

Philosophicus, IBZ — Internationale Bibliographie der Geistes- und

4304

Sozialwissenschaftlichen Zeitschriftenliteratur, Internationale Bibliographie der Rezensionen Geistes- und Sozialwissenschaftlicher Literatur, Linguistics

4305

and Language Behavior Abstracts, Mathematical Reviews, MathSciNet,

4306

Periodicals Contents Index, Philosopher's Index, Repertoire Bibliographique de la Philosophie, Russian Academy of Sciences Bibliographies.

4307

4308

4309

4310

## Contents

4311	GRAHAM STEVENS, <i>Metalinguistic Monstrosity and Displaced Communica-</i>	
4312	<i>tions</i> . . . . .	615
4313	RYAN COX, <i>Constitutivism about Instrumental Desire and Introspective Belief</i>	653
4314	ALIDA LIBERMAN, <i>The Mental States First Theory of Promising</i> . . . . .	683
4315	PETER FINOCCHIARO, <i>A Puzzle About Parsimony</i> . . . . .	711
4316	TOM DONALDSON, <i>David Armstrong on the Metaphysics of Mathematics</i> . . . .	729

PROOF