

# dialectica

International Journal of Philosophy

## Contents

JP SMIT & FILIP BUEKENS, <i>Is Somaliland a Country?: An Essay on Institutional Objects in the Social Sciences</i> . . . . .	1
LI ZHANG & LEON HORSTEN, <i>The Minimalist Theory of Truth and the Generalisation Problem</i> . . . . .	23
FR. JAMES DOMINIC ROONEY, OP, <i>The Problem of Thomistic Parts</i> . . . . .	45
WOLFGANG SPOHN, <i>A Generalization of the Reflection Principle</i> . . . . .	73
BATOUL HODROI, ANDREW JAMES LATHAM & KRISTIE MILLER, <i>Our Naïve Representation of Time and of the Open Future</i> . . . . .	97

# 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

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

## Review Editors

Stephan Leuenberger and Philipp Blum

## Editorial Committee

Sara Amighetti, Joshua Babic, Davood Bahjat, Philipp Blum (né Keller), Claudio Calosi, Alessandro Cecconi, Zoé Christoff, Fabrice Correia, Matthias Egg, Andrea Giananti, Martin Glazier, Aleks Knoks, Arturs Logins, Jörg Löschke, Giovanni Merlo, Robert Michels, Ryan Miller, Michael Müller, Paolo Natali, Donnchadh O'Conaill, Simone Olivadoti, Edgar Phillips, Stephanie Rennick, Sebastian Schmidt, Marco R. Schori, David Schroeren, Mike Stuart, Daniel Vanello.

## Consulting Board

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

March 2023

**Contents**

JP SMIT & FILIP BUEKENS, *Is Somaliland a Country?: An Essay on Institutional Objects in the Social Sciences* . . . . . 1

LI ZHANG & LEON HORSTEN, *The Minimalist Theory of Truth and the Generalisation Problem* . . . . . 23

FR. JAMES DOMINIC ROONEY, OP, *The Problem of Thomistic Parts* . . . . . 45

WOLFGANG SPOHN, *A Generalization of the Reflection Principle* . . . . . 73

BATOUL HODROJ, ANDREW JAMES LATHAM & KRISTIE MILLER, *Our Naïve Representation of Time and of the Open Future* . . . . . 97

PROOF

# Is Somaliland a Country?

## An Essay on Institutional Objects in the Social Sciences

JP SMIT & FILIP BUEKENS

Searle claims that his theory of institutional reality is particularly suitable as a theoretical scheme of individuation for work in the social sciences. We argue that this is not the case. The first problem with regulatory individuation is due to the familiar fact that institutional judgments have constrained revisability criteria. The second problem with regulatory individuation is due to the fact that institutions amend their declarative judgments based on the *inferential* (syntactic) properties of the judgments and in response to regulatory pressure, and not based on *descriptive* (semantic) properties and in response to matters of descriptive adequacy. These two problems imply that “regulatory kinds” (countries, borders, kings) will almost inevitably be disjunctive kinds that are ill-suited for scientific theorizing. This also explains why the law often makes odd pronouncements, e.g., calling ketchup a vegetable, considering an arm bent fifteen degrees to be straight, and not admitting that Somaliland is a country.

Somaliland is a democratically governed, autonomous region that maintains an independent police force, defends its borders, and issues currency in its own name.<sup>1</sup> Despite claims of statehood, it has not been officially recognized as a country by any state-level actors. Instead, it is considered an “autonomous region of Somalia.” Scotland is a semi-autonomous region that neither controls nor defends its borders<sup>2</sup> and does not govern its own affairs to the degree that countries typically do. Despite being atypical in these respects, it is officially recognized as a country.

---

1 See “Why Somaliland is not a recognized state” in *The Economist*, November 1, 2015.

2 Or, at least, the borders that *are* defended are defended *qua* United Kingdom and not *qua* Scotland.

25 Atypical cases like Somaliland, Scotland, and others immediately raise the  
 26 question as to the ontology of institutional objects like countries, presidents,  
 27 money, borders, and traffic lights. The issue is particularly pressing among  
 28 social scientists who study such phenomena. Suppose one is doing a cross-  
 29 country comparative analysis of some social or economic trend. The trend  
 30 does hold in Somalia (or the United Kingdom) but does not apply in Soma-  
 31 liland (or Scotland). In such a case, does Somaliland (or Scotland) serve as  
 32 counter-examples, thus weakening any potential claim to generality? Or does  
 33 Somaliland (or Scotland) not “count,” hence not affecting the generality of  
 34 any claim as to how widespread the trend actually is?<sup>3</sup>

35 The default answer one typically encounters when asking about what makes  
 36 it the case that *X* is a country is that *X* is a country if, and only if, regulative  
 37 bodies consider it a country. Of course, such regulative bodies have declared  
 38 that ketchup is a vegetable,<sup>4</sup> that Microsoft is a person,<sup>5</sup> and that an arm bent  
 39 15 degrees is straight.<sup>6</sup> Botany, psychology, and mathematics have ignored  
 40 these uses of “vegetable,” “person,” and “straight” to no ill effect. So why  
 41 should we care what regulative bodies have to say when individuating the  
 42 institutional world for the purposes of social science?

43 In this paper, we argue that social scientists should not feel compelled  
 44 to individuate the social world in the same way that institutions do. The  
 45 institutional use of language differs from the descriptive use proper to social  
 46 science in at least two ways, and both serve to make the regulatory schemes  
 47 of individuation used by institutions unsuited for descriptive work. The first  
 48 bad consequence of regulatory individuation is due to the familiar fact that  
 49 institutional judgments have constrained revisability criteria. This implies that  
 50 the facts picked out by institutional judgments will almost inevitably be non-  
 51 identical to the facts picked out by our best epistemic practices. The second

---

3 There are also more practical issues at stake. Somaliland, for instance, cannot receive state-level financial aid as such aid is earmarked for “countries” (Eubank 2015).

4 In 1981, the USDA’s Food and Nutrition Service recommended that schools could comply with official nutritional regulations by crediting condiments as vegetables. Although ketchup was not specifically mentioned (pickle relish was mentioned as an example), it became known as the “Ketchup as a vegetable” controversy.

5 The doctrine of corporate personhood grants entities like corporations some of the rights and obligations normally reserved for actual people.

6 The rule states that “[a] ball is fairly delivered in respect of the arm if, once the bowler’s arm has reached the level of the shoulder in the delivery swing, the elbow joint is not straightened partially or completely from that point until the ball has left the hand.” Yet an arm that does bend up to 15 degrees is not considered to violate this rule. The current laws are available at <https://www.lords.org/mcc/laws-of-cricket/>.

52 bad consequence of regulatory individuation is due to the fact that institutions  
 53 amend their declarative judgments based on the *syntactic* properties of the  
 54 judgments and in response to regulatory pressure, not based on *semantic*  
 55 properties and in response to matters of descriptive adequacy. This implies  
 56 that “regulatory kinds” (countries, borders, kings) will almost inevitably be  
 57 disjunctive kinds that are ill-suited for scientific theorizing. In making this  
 58 argument, we reject the account of Searle (2005), whose position implies that  
 59 social scientists should respect the individuation schemes of institutions.

## 60 **Searle on Institutional Facts**

61 John Searle, in a number of publications (1995, 2010), has defended an elegant  
 62 view of institutional facts. The object of study is institutional objects, i.e.,  
 63 objects that serve some social purpose in virtue of having certain deontic powers  
 64 (rights, duties, obligations). These deontic powers cannot be sufficiently  
 65 explained by the intrinsic or natural properties of the object itself, but are the  
 66 result of some institutional structure collectively endowing the object with  
 67 such properties by *recognizing* it as having such properties.

68 A typical example is that of a president. A president is not a president in  
 69 virtue of his or her physical or intrinsic properties, but in virtue of being  
 70 *recognized* as a president by the governing institution of the country that he  
 71 or she is president of. Paradigm cases of institutional objects also include  
 72 countries, borders, driver’s licenses, the playing field of a football game, and  
 73 so on. Our social reality is filled with such objects, and we interact with them  
 74 all the time.

75 Two aspects of Searle’s view are of particular interest. First, he claims  
 76 that institutional facts have the logical structure “*X counts as Y in C*” (1995,  
 77 28).<sup>7</sup> The *X*-term denotes the natural object, the *Y*-term is the institutional  
 78 specification of the object, and the *C*-term denotes the context in which  
 79 the institutional object has its function. In this way, Joe Biden counts as

7 A problem with this view is that the existence of some institutional facts does not seem to require the existence of anything for the *X*-term to denote. A paradigm case is money; most money does not exist in physical form but merely as account entries in bank ledgers. In response, Searle has stated that “*X counts as Y in C*” was only ever supposed to be a useful mnemonic that captures the core of his view (see Searle 2003). As “*X counts as Y in C*” is indeed a very useful mnemonic, and as nothing in the paper would be gained from using his later formulation, we stick with “*X counts as Y in C*.” (On the topic of the ontological status of money, see Smit, Buekens and du Plessis 2016, where we argue that Searle’s *X*-term can be interpreted as referring to an abstract object.)

80 the president of the United States at present, a specific line counts as the  
 81 goal line during a game of football, and so on. Second, Searle claims that  
 82 the recognition that is constitutive of the existence of an institutional fact is  
 83 essentially *collective* recognition. Institutions are collectives, and the collective  
 84 recognition by which they endow an object with deontic powers is irreducible<sup>8</sup>  
 85 to individual recognition (1995, 24–25).<sup>9</sup>

86 Searle (2005) considers his view to have particular relevance for the social  
 87 sciences. He introduces an article on the relevance of his view for economics  
 88 (and social science in general) as follows:

89 When I was an undergraduate at Oxford, we were taught eco-  
 90 nomics almost as though it were a natural science. The subject  
 91 matter of economics may be different from physics, but only in  
 92 the way that the subject matter of chemistry or biology is different  
 93 from physics. The actual results were presented to us as if they  
 94 were scientific theories. So, when we learned that savings equals  
 95 investment, it was taught in the same tone of voice as one teaches  
 96 that force equals mass times acceleration. And we learned that  
 97 rational entrepreneurs sell where marginal cost equals marginal  
 98 revenue in the way that we once learned that bodies attract in  
 99 a way that is directly proportional to their mass and inversely  
 100 proportional to the square of the distance between them. At no  
 101 point was it ever suggested that the reality described by economic  
 102 theory was dependent on human beliefs and other attitudes in a  
 103 way that was totally unlike the reality described by physics and  
 104 chemistry. (Searle 2005, 1)

105 Searle sets up a basic distinction between the objects of the physical sciences  
 106 and the objects of social sciences and advises social scientists to heed the  
 107 fact that their objects are fundamentally unlike those of the physical sci-  
 108 ences. The objects of social science are frequently institutional objects and, as

---

8 For a critique of this claim, see Smit, Buekens and du Plessis (2011, 2014), where we develop the incentive account of institutional facts. On our view, institutions can be fully understood in terms of incentives and actions, and the recognition of such incentives and actions need not be collective. The view is similar to Guala and Hindriks (2015; Hindriks and Guala 2015)—also see Guala (2016)—who account for institutions in terms of rules in game theoretical equilibria.

9 Searle, in recent years, has recognized that, in some cases, forms of collective institutional recognition may reduce to individual recognition (2010, 58).



109 such, should be understood as explained above, i.e., in terms of the collective  
110 recognition of objects as having certain deontic powers.

111 Of particular importance to the current discussion is his claim that such  
112 objects can only exist for as long as they are represented as existing (1995,  
113 13) and his claim that such objects should be understood as having a logical  
114 structure (1995, 22). This implies that if one is a social scientist and wishes to  
115 study borders, countries, or presidents, then one must understand one's area  
116 of study as pertaining to those things *recognized* to be borders, countries, and  
117 presidents. In other words, those things that exist in virtue of the collective  
118 acceptance of a declaration of the form "X counts as Y in C." This implies,  
119 although it is not explicitly stated by Searle, that the social scientist must  
120 individuate the institutional world as the institutions that create it do. For  
121 if it is constitutive of borders, countries, and presidents that they must be  
122 *recognized* to be these objects, then studying objects not so recognized is to  
123 not study borders, countries, and presidents at all.

## 122 Two Peculiarities of the Institutional Use of Language

### 2251 A Toy Example

126 Below, we will argue that social scientists would be ill-served if they employed  
127 the individuation schemes used by institutions themselves. The argument is  
128 based on the fact that institutions use language in peculiar ways, and these  
129 ways make institutional standards of individuation ill-suited for the purposes  
130 of scientific description. This is not to say that there is any specific problem  
131 with describing institutions; rather, the claim is that the social scientist should  
132 not feel compelled to use the regulatory schemes of individuation adopted by  
133 institutions when describing institutional facts.

134 For purposes of exposition and illustration, it will be useful to have a toy  
135 example at our disposal. Suppose there is a village in the Scottish Highlands  
136 that has a cultural ritual called "Firecasting" that takes place annually on the  
137 first day of spring. They celebrate the end of winter and the reduced need  
138 for heating by letting each member of the village attempt to light a torch on  
139 fire, run to the firecasting line, and hurl it into a lake, extinguishing the flame  
140 within twenty seconds. Those who succeed get a medal (and receiving such a  
141 medal has significant prestige in the village).

142 In Firecasting, there is an umpire who keeps time, adjudicates whether a  
 143 flame has been extinguished, etc. Every time a torch has been extinguished, the  
 144 umpire proclaims, “Player *x* is a firecaster,” i.e., a flame has been extinguished.

## 2.4.2 *First Peculiarity Constrained Revisability*

146 The umpire in Firecasting has to judge whether a specific state of affairs  
 147 obtains, namely, whether the flame has been extinguished. This is a judgment  
 148 that any spectator can also make. The umpire’s judgment, however, counts  
 149 in a way that the judgments of spectators do not. Consider the following  
 150 judgment:

151 (1) John is a firecaster.

152 If a spectator makes a judgment by using (1), then this is a speech act of  
 153 description that asserts that John extinguished the flaming torch. As it is a  
 154 standard instance of description, it can be straightforwardly true or false. Call  
 155 the use of the institutional term “firecaster” in a speech act of description the  
 156 *descriptive use* of the term.

157 If (1) is used by the umpire, however, the situation is different. The um-  
 158 pire’s use of (1) is *based on* his assessment of whether the flame has been  
 159 extinguished, yet his speech act is that of declaration. His speech act has the  
 160 function of creating a certain institutional fact, namely, the institutional fact  
 161 that John is a firecaster. In Searle’s terminology, such a judgment has certain  
 162 deontic consequences, namely, that the player is entitled to be awarded a  
 163 medal by the village. Call such a use of the institutional term “firecasted” the  
 164 *regulative use*.<sup>10</sup>

165 Note that if the umpire makes a mistake in adjudging whether John has  
 166 extinguished the flame and erroneously declares that he is a firecaster, then  
 167 the descriptive content of (1) is false, yet the regulative content of (1) can still  
 168 be affirmed. This is so that, even if the umpire makes a mistake, the deontic  
 169 consequences of his judgment will still obtain, i.e., John will still be entitled  
 170 to the medal. What the umpire commits himself (and the village) to through  
 171 the speech act of declaration is, above all else, that John is entitled to receive  
 172 the medal. The umpire’s judgment might be *based on* whether the descriptive

---

10 Of course, (1) can also be used in a third way: as a *report* of an umpiring judgment. This use, while also descriptive, is distinct from the descriptive use in the main text and need not trouble us here.

173 content of (1) obtains, yet what is affirmed by the umpire in making his  
174 judgment is something else.<sup>11</sup>

175 It is a staple of the literature on the philosophy of law (see, for instance,  
176 Hart 1961) that institutional declarations cannot be revised in light of future  
177 evidence in the same way that descriptive judgments can. Even if the umpire  
178 and the village see conclusive evidence that John did not extinguish the torch,  
179 they may choose to “let the judgment stand,” i.e., to remain committed to  
180 enforcing the deontic consequences of the original regulative judgment. While  
181 the village may choose to explicitly adopt regulative rules that do allow for the  
182 revision of prior judgments, there is nothing inherently irrational about not  
183 doing so, as affirming the regulative content of (1) does not logically commit  
184 them to any specific position as to the truth value of the descriptive content of  
185 (1). In this way, the regulative judgment contrasts sharply with the descriptive  
186 judgment, as it is a *sine qua non* of descriptive practice that such judgments  
187 are always revisable in light of future evidence.

188 In fact, most real-world sports do not, except in extreme cases, allow such  
189 later reviews of umpiring decisions. Mistakes inevitably happen, and sports  
190 fans everywhere use the institutional term itself in a descriptive way in order  
191 to register their disagreement with the referee. Consider judgments like “That  
192 was never a strike!” “He was miles offside!” “It pitched outside leg stump!”  
193 and so on. In such cases, the utterer uses a *non-institutional*, descriptive  
194 standard for applying the terms “strike,” “offside,” and “outside.” Here, the  
195 institutional term is used in order to voice disagreement with the factual basis  
196 of an umpiring decision (and also to draw attention to the unfairness of the  
197 deontic consequences of such a decision).

198 The phenomenon of *constrained revisability* is found in all institutional  
199 settings. While institutional judgments can sometimes be overruled—i.e.,  
200 appealed in various ways—such revisability is constrained in a way that that  
201 open-ended, epistemic inquiry is not. For example, legal systems in a wide  
202 variety of countries recognize a principle of “double jeopardy,” whereby an  
203 accused cannot be retried for an offense that they have already been acquitted  
204 of. This remains so even if definitive evidence of prior guilt is produced and

---

11 This distinction between the basis for an institutional judgment and its deontic consequences was first set out in Ransdell (1971, 388). I am departing from his terminology (he distinguishes between the “connotation” of a term and its “import”), but this departure should not be taken to imply any difference of substance. Ásta draws a similar distinction between “base properties” and “conferred properties” (2018).

no one believes that the descriptive judgment underlying the institutional declaration was accurate.

Some legal systems do allow various, tightly restricted exceptions to this principle. In general, though, the revisability of the legal declaration is constrained in a way that commitment to the underlying descriptive claim is not. The distinction between the descriptive use and the regulative use of legal terms is again well-recognized in our ordinary discourse. Consider judgments like “Andy Dufresne was innocent,”<sup>12</sup> “OJ Simpson was guilty,” “Jimmy Saville was a criminal,” and so on.<sup>13</sup>

The fact that (1) can express distinct speech acts with distinct criteria of revisability means that the denotation of the descriptive use and the regulative use of a term can diverge. In the case of the game Firecasting, the denotation of the descriptive use of the term “firecasters” will include those who succeeded in extinguishing a flame. The denotation of the regulative use will include all those who were *adjudged* to have extinguished a flame. If a scientist were to study firecasting, then the nature of the study might force her to take the distinction seriously and use one or the other criterion. In this way, if the scientist were tasked with determining what physical characteristics allow one to firecast, then she would be ill-served by the regulative use. This is because, if a number of serious umpiring errors have occurred in the history of Firecasting, any law-like generalization that the scientist seeks to uncover will be much more likely to apply to those who *actually* achieved the feat of extinguishing a flame and not merely those adjudged to have done so. The denotation of the regulative use of “firecasting” will almost inevitably be non-identical to the denotation of the descriptive use of “firecasting”; the denotation of the former will be more heterogeneous with regards to physical characteristics (as it includes both those who succeeded and those who did not) and, as such, less likely to be the object of useful law-like generalizations of the required type.

The opposite is likely to be true for the historian of the game. The historian who writes about the stars of the game is implicitly, and correctly, writing about the regulative use when writing about those *recognized* to have fire-

12 The protagonist of the Stephen King novella *Rita Hayworth and the Shawshank Redemption*—later made into the film *The Shawshank Redemption*—who was convicted of a crime he did not commit.

13 The disagreement need not take the form of a factual disagreement but can also be used to express disagreement with the normative judgment behind an institutional judgment. Few people would consider Nelson Mandela a “terrorist,” despite the fact that he used to be on the US terrorist watch list.

237 casted. Here, the main interest lies in those falling under the denotation of  
238 the regulative use, and as such, reports of prior regulative use are appropriate  
239 to the study. As this is the main topic of interest, any law-like regularities that  
240 the typical sports historian seeks are supposed to concern those who were  
241 *adjudged* to have firecasted.

242 The same distinction applies to the academic study of less frivolous matters.  
243 Consider a criminologist who aims to make discoveries about the causes of  
244 crime in order to determine how law-breaking can be prevented. Here, the  
245 interest is likely to lie in determining what causes people to break the law.  
246 Breaking the law, of course, is not the same thing as being *adjudged* to have  
247 broken the law. In this manner, the criminologist would feel vindicated if their  
248 theory applied to someone who did commit a crime but was never caught  
249 or convicted. In the same way, they would be untroubled if their theory did  
250 not apply to someone who was wrongfully convicted. When we express an  
251 interest in preventing crime, we are typically not expressing an interest in  
252 having fewer of the people who break the law caught but in having the law  
253 broken less. This makes the *descriptive use* of the term “criminal” the one  
254 appropriate to such a study.

255 The criminologist could, of course, decide to try and find out what distin-  
256 guishes those convicted of committing a crime from those not so convicted.  
257 Here, the regulative use of “criminal” would be appropriate to the study. Note,  
258 however, that we have some reason to believe that in typical cases, law-like  
259 generalizations are more likely to apply to the descriptive use. The denotation  
260 of the descriptive use of criminal would include all those who broke the law.  
261 The denotation of the regulative use would include those who broke the law  
262 and were convicted and those who did not break the law and were convicted,  
263 while excluding those who did break the law and were not convicted. As  
264 the latter category is individuated in terms of a more heterogeneous mix of  
265 properties, one would suspect, *ceteris paribus*, that the descriptive use of the  
266 term “criminal” would be more suitable for obtaining law-like generaliza-  
267 tion. Simply put, it will typically be easier to obtain general truths among a  
268 group whose members were rightfully classified as belonging to the group  
269 than among a group that includes a mixture of those correctly and incorrectly  
270 classified as members of the relevant group.

271 The above reasoning implies that the constrained revisability of regulative  
272 judgments sometimes gives the social scientist a good reason to, despite using  
273 the terms used by some specific institution itself, reject the individuation  
274 scheme of the institution. This is so as, firstly, the ultimate goal of the inquiry

275 (i.e., crime prevention) may demand it. Secondly, the descriptive use of the  
 276 institutional term will be more suitable for law-like generalization and, so,  
 277 more useful to social science.

283 *Second Peculiarity Institutional Judgments Are Amended Based on*  
 279 *Their Inferential (Syntactic) Properties, Not Descriptive (Semantic)*  
 280 *Properties*

281 Suppose that the village that practices Firecasting notices that players are  
 282 sometimes prevented from hurling the torch by other players kicking them  
 283 just as they are about to hurl it and, in so doing, making it less likely that  
 284 the player throws across the line. They wish to make such behavior pointless  
 285 and, so, announce that players who are kicked as they are about to throw will  
 286 receive a medal anyway, even if their throw did not cross the line.

287 The required rule change can be made in two distinct ways. Prior to the  
 288 rule change, the relevant rules of Firecasting are as follows:

- 289 (1) A player  $x$  has firecasted if, and only if,  $x$  is adjudged to have extin-  
 290 guished the flaming torch by hurling it into the ocean within the context  
 291 of the game of Firecasting.  
 292 (2) Firecasters are entitled to receive a medal from the village.

293 The first way to change the rule so as to award those who were kicked just  
 294 prior to throwing would be to amend the definition of firecasting so that  
 295 those who were kicked also “count” as firecasters. This option is analogous  
 296 to a “penalty try” in rugby. If a rugby player is illegitimately prevented from  
 297 scoring a regular try, the referee may award a so-called “penalty try” to the  
 298 team prevented from scoring. A rugby team awarded a penalty try is awarded  
 299 five points in the same way that a team that scores a regular try is awarded  
 300 five points. In this way, the penalty try “still counts,” despite the fact that the  
 301 attacking team was prevented from scoring a regular try. In the same spirit,  
 302 the village can amend (1) as follows:

- 303 (1\*) A player  $x$  has firecasted if, and only if,  $x$  is adjudged to have extin-  
 304 guished the flaming torch by hurling it into the ocean or  $x$  is adjudged  
 305 to have been kicked prior to hurling the flaming torch within the context  
 306 of a game of firecasting.

307 The village, however, need not amend the definition of “firecasting” in order  
 308 to secure the result that those who are kicked in order to prevent them from  
 309 firecasting still receive a medal. They can leave (1) intact and simply amend  
 310 (2) so that it states that those who are kicked also receive a medal from the  
 311 village. In this way, (2) can be amended as follows:

312 (2\*) Firecasters and those who were adjudged to have been kicked just prior  
 313 to hurling the flaming torch within the context of a game of Firecasting  
 314 are entitled to receive a medal from the village.

315 (2\*) directly regulates the result of the person attempting to firecast being  
 316 kicked, whereas (1\*) does indirectly by changing the concept of “firecasting.”  
 317 Yet these two ways of amending the rules are equivalent; both of the above  
 318 rule changes would have the effect that those who were kicked receive a medal.  
 319 The change can be made in distinct ways, as, in the context of enforcement  
 320 of such rules, the rules of Firecasting constitute a set of *premises* that can be  
 321 amended in distinct ways so as to, in conjunction with judgments about an  
 322 instance of the game, imply the statement that some player who has been  
 323 kicked should receive a medal. In other words, the aim of the village, when  
 324 amending the rules of Firecasting, is to appropriately link the following two  
 325 statements concerning some specific instance of the game.

326 (3) Player *K* was adjudged to have been kicked prior to hurling the flaming  
 327 torch within the context of a game of Firecasting.

328 (4) Player *K* is entitled to receive a medal from the village.

329 (3) is a specific judgment concerning some specific instance of the game of  
 330 Firecasting, and (4) is the regulatory response to what was adjudged to have  
 331 happened in some such specific instance of the game. The village aims to  
 332 formulate rules that, in conjunction with (3), imply (4). The combination of  
 333 (1\*), (2), and (3) implies (4), and the combination of (2\*) and (3) also implies  
 334 (4). In this way, the fact that the rules can be amended in distinct ways reflects  
 335 no more than the fact that the same conclusion can follow from distinct sets  
 336 of premises.

337 In the above case, the first way of changing the rules amounted to changing  
 338 the definition of “firecasting,” whereas the second amounted to changing the  
 339 statement of rewards given out by the village. The regulatory equivalence  
 340 of these changes in our toy example is a phenomenon that applies to law  
 341 in general. When we wish to amend the law in order to secure a specific

342 consequence, there will always be distinct ways of doing so, and the only  
343 criteria for choosing whether to amend a definition or amend some statement  
344 of penalties or awards is, where rational, pragmatic.

345 Legal language turns out to be holistic in an almost Quinean way (Quine  
346 1951). The law is holistic in two distinct ways. First, there is no one correct  
347 way to change the law so as to secure some regulatory response. Second, the  
348 list of claims we call definitions has no special status that prevents them from  
349 being changed so as to secure the desired regulatory response.

350 The fact that the law is holistic serves to explain why the law often uses per-  
351 fectly ordinary terms in peculiar ways. The claim that a company is a “person”  
352 is just a tool to secure a regulatory response concerning the legal liability of  
353 the members of a corporation; the claim that ketchup is a “vegetable” is a tool  
354 to effectively lower the legally mandated nutritional requirements for school  
355 lunches. In the same way, the claim that an arm bent 15 degrees is “straight”  
356 is a tool to secure the result that cricketers may bowl with a slightly bent arm.

357 Cases of atypical use of terms like “ketchup,” “person,” “straight,” and  
358 the like serve to demonstrate something important. When the lawmaker  
359 changes the law, it has no overriding reason to respect the semantics of the  
360 term (as used in non-legal contexts). Rather, specific statements only matter  
361 inasmuch as they help to, in conjunction with other statements, secure the  
362 desired regulatory response when the law is applied. Such regulatory responses  
363 in specific instances can be represented as the conclusions of arguments  
364 that have legal statements among their premises. As the overriding factor  
365 governing the formulation and emendation of laws is the regulatory response  
366 to which it gives rise, this implies that the overriding factor governing the  
367 emendation of statements within a system of law is *the role of such statements*  
368 *in facilitating inference*. This, in turn, implies that we can expect changes  
369 to the law to end up radically changing the denotation of terms that were  
370 originally used in a perfectly ordinary sense. In the final analysis, this is due to  
371 the fact that institutional judgments are amended in virtue of their inferential  
372 (*syntactic*) properties and not their descriptive (*semantic*) properties.

373 The first peculiarity of language that was noted was that constrained re-  
374 visability meant that the denotation of the regulative use of a term would  
375 not exactly coincide with the denotation of the descriptive condition based  
376 on which the term is applied. The second problem, however, is much more  
377 basic and would apply even if judges never made mistakes. Lawmakers will  
378 change the content of perfectly ordinary terms in order to secure regulatory  
379 consequences. This implies that the legal system will tend towards a scheme



380 of individuation designed to serve regulatory and not descriptive purposes.  
381 This much is obvious enough, but it has the less commonly understood conse-  
382 quence that institutional judgments will be amended based on their syntactic  
383 properties and not on their semantic properties.

384 This implies that technical terms introduced for some regulative purpose  
385 (like “firecasting”) are not constrained so as to include only relevantly similar  
386 elements under their denotation. Furthermore, even when the regulative term  
387 is taken from ordinary language (like “straight,” “vegetable,” etc.), the term will  
388 tend to start out being used as legal terms in their familiar sense but will often  
389 end up including unlike objects in the same category. In this way, “regulatory  
390 kinds” will end up, if judged against a standard of descriptive adequacy,  
391 becoming disjunctive kinds that are ill-suited to scientific theorizing.

392 No one would expect mathematicians to do useful work while treating an  
393 arm bent 15 degrees as straight, and no one would expect a botanist to employ  
394 the term “vegetable” so as to include ketchup. These cases, however, are just  
395 the tip of the iceberg that serve to make the general phenomenon visible. The  
396 physical sciences pay no attention to regulative bodies when individuating the  
397 world, as such bodies are simply involved in another kind of activity altogether.  
398 In the same way, there is no reason for the social scientist to consider herself  
399 uniquely encumbered by, and beholden to, a schema of individuation that  
400 does not serve her purposes.

401 Note that the point concerning ordinary terms being introduced into law is  
402 not merely that such regulatory terms “change their meaning.” The problem,  
403 rather, is that such changes occur due to inferential (syntactic) considerations.  
404 Legal changes may be phrased as changes in definition (as when the definition  
405 of “firecasting” is amended so as to include being kicked) or as changes in  
406 regulation (such as when those kicked during firecasting also receive a reward).  
407 Whether these changes are phrased as one or the other change has little, if  
408 anything, to do with descriptive adequacy, and so we end up with categories  
409 that mix unlike things together, i.e., “regulatory kinds” become disjunctive  
410 kinds.

411 To illustrate the above point, consider the difference between the descriptive  
412 term “computer” and the regulative term “king.” The term “computer” was  
413 originally used to denote people, specifically those employed to engage in  
414 tedious tasks of rote calculation.<sup>14</sup> The project to mechanize such tasks was  
415 originally described as the project of creating a “mechanical computer,” and

---

14 For an interesting history of pre-mechanical calculation, see Grier (2005).

416 this description was no mere tautology. Once the project succeeded, however,  
417 and the human computers disappeared, the meaning of the term “computer”  
418 underwent a social shift until it denoted only machines designed to perform  
419 such calculations. In fact, the change in the use of the term exhibits a nice  
420 symmetry; today, if we call someone a “computer,” it is a metaphorical use of  
421 the term that suggests extreme proficiency at calculation, a tendency to act  
422 without emotion, or some such.

423 The above change is generally socially recognized as a fundamental change  
424 in the *meaning* of the term “computer.” In principle, we could have treated  
425 the shift differently, for instance, by saying that the term merely expanded  
426 its denotation so as to include both human and mechanical, and eventually  
427 electronic, computers. There would be no point in doing so, however, for then  
428 “computer” would be a disjunctive kind that groups two radically distinct  
429 kinds of things together. Our descriptive language is guided by descriptive  
430 adequacy, and so simply treating the content of the term as having changed  
431 completely individuates the world in a much more useful way.

432 The same is not true for the institutional pair “king”/“queen” when used in  
433 a regulative manner. Kings and queens, historically, are paradigmatic persons  
434 who rule a state by decree and who obtain their position by right of birth.  
435 Today, however, in the vast majority of countries that still recognize a “king”  
436 or “queen,” being a king or queen is primarily a symbolic or ceremonial role.  
437 While today’s kings and queens do have some influence, this influence is so  
438 different in kind from the right to rule by decree that the two kinds of “king”  
439 or “queen” are beyond any meaningful similarity or comparison. While we  
440 may loosely say that “the meaning of being a king or queen has changed,”  
441 we do not generally consider the term to have changed its semantic content  
442 in the same way that the term “computer” has. This is despite the fact that  
443 the term categorizes together entities with vastly different social roles. If a  
444 historian or sociologist were to uncritically accept the institutional use of the  
445 term “king” (or “queen”) and try to determine commonalities or differences  
446 between kings, the very category of analysis would serve to unnecessarily  
447 complicate the inquiry. The term would group together those who ruled by  
448 decree as well as those whose social role is effectively a more dignified version  
449 of a mascot. If our purpose is descriptive adequacy, then little is to be gained  
450 by an individuation scheme that treats those who ruled by decree (old-style  
451 kings, the present-day King of Swaziland, etc.) with the current Queen of  
452 England or the current Queen of Denmark. Furthermore, it would exclude  
453 those whose social role is similar to that of old-style kings and queens, i.e.,

454 dictators who de facto rule by decree and have their position in virtue of birth,  
 455 e.g., the North Korean leader Kim Jong-Un.

456 The point of the above is not to criticize the existence of present-day royalty  
 457 or to suggest a change in linguistic habit. The point, rather, is that a social  
 458 scientist who accepts an individuation scheme in which an all-powerful king  
 459 and the current Queen of Denmark are the same sort of thing<sup>15</sup> is as absurd  
 460 as a mathematician who treats all lines that bend less than 15 degrees as  
 461 “straight.”

462 The above considerations concern both cross-institutional identification  
 463 (i.e., whether the current King of Denmark and the current King of Swaziland  
 464 are the same sort of thing) and inter-temporal identification (i.e., whether  
 465 the kings and queens of centuries ago are the same sorts of things as the  
 466 current Queen of England). It is a fundamental constraint upon inquiry that  
 467 our criteria of individuation remain *constant*, and here, the regulative use  
 468 of institutional terms is a poor guide to scientific individuating practices. It  
 469 is for this reason that the social scientist should feel under no obligation to  
 470 accept institutional standards of individuation; in fact, she should rearrange  
 471 the conceptual world as she sees fit.<sup>16</sup>

### 472 **Cases Where the Regulative Use of Institutional Terms Is** 473 **the Correct Use**

474 The point of the above is not that social scientists should never employ the  
 475 regulative use of institutional terms as basic terms of inquiry. Institutions  
 476 do manage to affect the world via declarations, and Searle is correct that this  
 477 phenomenon is important to understanding our social and institutional world.  
 478 We can distinguish three reasons for adopting the regulative use as a term of  
 479 inquiry.

480 First, our interest may lie precisely in *the objects grouped together* by the  
 481 regulative use of an institutional term. In this way, as mentioned earlier, we  
 482 may wish to inquire into the difference between those who are convicted of

---

15 One could object that kings and queens do form a kind in virtue of their genetic relation to an ancestor. This is so and means that the term, so construed, would be useful for geneticists. Most of the time, however, when considering kings and queens, our interest lies in their social role, and here, the regulative use of the term is a plain obstacle to inquiry.

16 Our account has the additional advantage that it does not overemphasize the role of normativity in the causal processes operative in social reality. See Turner (2010) and Guala (2015) for critical assessments of the (over)use of normativity in the social sciences and social ontology.

483 committing a crime and those who, while having committed a crime, are  
 484 acquitted. Or, alternatively, we may be interested in the difference in severity  
 485 of sentence among those convicted of a crime. Such topics are a staple of  
 486 criminological and sociological studies that try and determine what effect  
 487 categories of identity (race, gender, etc.) or socioeconomic attainment have  
 488 on rates of conviction and severity of sentence. In such cases, our interest lies  
 489 precisely in a category that exists in virtue of regulative declarations; hence,  
 490 the regulative category is proper to the study.<sup>17</sup>

491 Second, the declarations made by institutions have a *causal impact* on  
 492 the world, and our interest may lie precisely in studying the effect of such  
 493 an impact. In this way, the criminologist may be interested precisely in the  
 494 impact of a criminal conviction on one's life prospects. In this case, again,  
 495 the regulative use is proper to the study in virtue of the causal role of such  
 496 regulative judgments.

497 An interesting sub-class of the causal impact that institutional declarations  
 498 can have is where such declarations have a *symbolic impact* on the objects  
 499 of such a declaration; we may well be interested in studying the nature and  
 500 effects of such a symbolic impact. In this way, a historian or sociologist may  
 501 be interested in changes in self-conception that occur among those people  
 502 occupying a territory that is widely recognized as being a "country," or changes  
 503 in self-conception among those recognized as "criminals," and so on.<sup>18</sup>

504 Third, the social scientist may be interested in a category that need not be  
 505 governed by regulative use, but the institutional use is close enough, for the  
 506 purposes of the study, to what they are trying to identify that it is a *useful proxy*  
 507 for the descriptive use. If a country's rules remain relatively stable over time,  
 508 the judiciary does a decent job of applying the laws, and the concepts involved  
 509 happen to individuate reality in a descriptively useful way, the institutional  
 510 category should be good enough for useful inquiry. Consider, for instance,  
 511 a scientist who wishes to study whether the color of a motor vehicle has an  
 512 impact on people's propensity to speed. Strictly speaking, some people who  
 513 speed will not be among those convicted of speeding, whereas some of those

---

17 See Wilson (2007) for a related argument that the importance of Searle's work to the social sciences is more limited than one might suppose.

18 In this paper, we mostly speak of classification as a matter of putting objects with similar causal powers together. An anonymous referee points out that the social sciences do more than try to arrive at law-like generalizations. Nothing in our argument prohibits non-causal schemes of individuation that may prove useful in interpretive or normative projects; the point is that the Searlean project does not tie our hands.

514 convicted will have been innocent. But if all the scientist is looking for is  
515 a rough correlation in aggregate and the legal system has been reasonably  
516 efficient, then counting all those convicted of speeding as “speeders” should  
517 be a good enough sample to do meaningful statistical work.

518 In endorsing the above regulative uses, we also embrace something close  
519 to pluralism about general institutional terms. Good usage will be polysemic;  
520 the social scientist will inevitably have to craft the terms of their inquiry to  
521 the topic at hand. What we object to, on the grounds discussed, is the idea  
522 that regulative bodies should be implicitly granted the power to set the terms  
523 of our descriptive agenda.<sup>19</sup>

#### 524 **4 Is Somaliland a Country?**

525 Somaliland is a country. More specifically, we think that, except in the very  
526 specific types of cases previously explained, i.e., cases where our epistemic  
527 interest is precisely in those objects grouped together by institutional declara-  
528 tions, the social scientist should view Somaliland as a country. Here, we do  
529 not base this claim on any specific definition of the term “country.” Rather,  
530 our judgment reflects the fact that Somaliland, once we ignore regulatory  
531 schemes of individuation for the reasons outlined in this paper, seems entirely  
532 like paradigm cases of countries, i.e., Kenya, Germany, Chile, Japan, etc.

533 In this paper, we have explained why we think that the social scientist  
534 should, in principle, be very wary of adopting institutional schemes of individ-  
535 uation. This matters, as the currently dominant theory of institutions, i.e., the  
536 Searlean theory, effectively adopts and legitimizes institutional schemes of  
537 individuation, and, hence, it is worth knowing why the social scientist should  
538 feel free to disregard Searle’s view. This is so, especially as Searle explicitly  
539 recommends that social science adopt his theory of institutions (Searle 2005).  
540 Also note that *reflexive* definitions, i.e., definitions on which an entity gains its  
541 identity from being *considered* to be the things that it is, long predate Searle.

542 More important, however, is the question as to the scope of the problem, i.e.,  
543 the question of how much harm is done by social scientists adopting regulatory  
544 schemes of individuation. Our argument is compatible with *quietism* about  
545 regulatory individuation, i.e., the view that Somaliland is an edge case, a mere  
546 curiosity whose exclusion from the list of countries does no real harm to  
547 social analysis. Our view is also compatible with *revisionism* about regulatory

---

19 We would like to thank an anonymous referee for pressing us to be explicit on this point.

548 individuation, i.e., the view that Somaliland and cases like it serve to make  
 549 visible a deep problem that calls for social scientists to abandon regulatory  
 550 schemes of individuation in favor of a series of successor concepts more suited  
 551 to descriptive purposes.

552 The question of which position on the continuum between quietism and  
 553 revisionism is most justifiable is beyond the ambition of the present work.  
 554 We can see the appeal of quietism; it would appear ridiculous to expect an  
 555 economist writing about the correlation between countries in the measured  
 556 link between inflation and unemployment to worry too much about whether  
 557 his fundamental categories of analysis are making his job harder than it needs  
 558 to be.

559 We can also, however, see the appeal of revisionism. Consider the definition  
 560 below, intended to capture the regulatory notion of a “country”:<sup>20</sup>

561 A country is a region that is identified as a distinct entity in political  
 562 geography. A country may be an independent sovereign state or part  
 563 of a larger state, as a non-sovereign or formerly sovereign political  
 564 division, or a geographic region associated with sets of previously  
 565 independent or differently associated people with distinct political  
 566 characteristics.

567 The above definition—in addition to being vague—is disjunctive to an extreme  
 568 degree. It is akin to a definition of “vegetable” that includes not only ketchup  
 569 but also all bottles of Worcestershire sauce that are older than three months.  
 570 The problem with such a disjunction is plain: What possible reason could we  
 571 have to expect that some underlying causal process could produce similar  
 572 effects across entities that have been grouped together merely as a matter of a  
 573 series of historical regulatory contingencies?

574 Current practice seems to imply at least some deviation from quietism.  
 575 Social scientists are not naive and have not stayed slavishly faithful to institu-  
 576 tional categories of individuation. The *CIA World Factbook*, for example, on  
 577 its list of countries by gross domestic product, does not list England, Scotland,  
 578 or Wales among the entries, even though they are generally recognized as  
 579 countries.<sup>21</sup> It does, however, list the European Union, despite the fact that it  
 580 is not recognized as being a country. This makes sense, as the interest of the

---

20 From [worlddata.info](http://worlddata.info).

21 *CIA World Factbook* is available at: <https://www.cia.gov/library/publications/resources/the-world-factbook/rankorder/2001rank.html>.

581 economist would be in finding a category of individuation that identifies indi-  
 582 vidual units of action, i.e., units with a fair degree of autonomy *qua* matters  
 583 of economic production and exchange. When it comes to such practices, the  
 584 present paper serves to justify how such deviations from institutional schemes  
 585 of individuation are, *contra* Searle, perfectly justified.<sup>22</sup>

586 The question, however, is whether current practice occupies the appropriate  
 587 position on the continuum between revisionism and quietism. Note that the  
 588 *CIA World Factbook* does not list Somaliland, despite there being very little  
 589 reason not to do so once we abandon the purely regulatory use of “country.”  
 590 In fact, once we take the matter of individuation seriously, we may well  
 591 have reason to include sub-units of various “countries” under the de facto  
 592 control of some entity other than the recognized government, i.e., parts of  
 593 “countries” under the control of rebel groups or drug cartels.<sup>23</sup> This may sound  
 594 radical, but if our interest lies in discovering the units of political and/or  
 595 economic action—and hence in groups that have a high degree of autonomy  
 596 over running their own affairs—then there is little reason to exclude them.  
 597 We may well learn interesting things by considering such entities *qua* units  
 598 of economic and political action, for they are effectively no different from  
 599 “countries” under military or dictatorial control.

## 605 5 Conclusion

601 In this paper, we have argued that social scientists should not be Searleans  
 602 when it comes to their own categories of analysis, i.e., be wary of employing  
 603 the regulative use of institutional terms for purposes of individuation. There  
 604 are two main problems. First, the revisability of institutional judgments is non-  
 605 epistemically constrained, i.e., mistakes do not get corrected in the same way  
 606 that we correct them when dealing with descriptive assertions. This means  
 607 that the social scientist would frequently be better served by employing the  
 608 descriptive use, and not the regulative use, of institutional terms as a basis of

22 The Searlean could respond by saying that such usage of “country” is a mere loose usage, done for practical purposes. Such a response, however, opens up the line of attack that we have been pressing, for it implicitly admits that the Searlean scheme of individuation is not suited to social science. We would like to thank an anonymous referee for pressing us on this point.

23 The CIA estimates that roughly 20% of Mexico is under the control of the drug cartels. (See “Mexico’s government control threatened by criminal groups claiming more territory” in *The Washington Post*, October 29, 2020.) Interestingly, some such drug cartels engage in activities commonly associated with governments, e.g., the provision of social services. See Flanigan (2014) for a discussion of this phenomenon.

609 individuation. The second problem, and by far the most important one, is due  
 610 to the fact that institutions individuate in order to regulate, not to describe.  
 611 Such regulation is holistic, and, hence, the usage of terms will change based on  
 612 their *syntactic* properties and in response to regulatory pressure, not based on  
 613 their *semantic* properties and in response to matters of descriptive adequacy.<sup>24</sup>  
 614 This means that “regulative kinds” will inevitably tend to become disjunctive  
 615 kinds, and so the law will be prone to the seeming absurdity of classifying  
 616 ketchup as a vegetable, considering an arm bent 15 degrees to be straight, and  
 617 so on. This implies that the social scientist will sometimes be better advised  
 618 to ignore *both* the descriptive and regulative use of institutional terms and  
 619 to invent institutional categories that have never been subject to regulative  
 620 declaration at all.\*

621 JP Smit

622  0000-0003-1524-8845

623 University of Stellenbosch

624 jps@sun.ac.za

625 Filip Buekens

626  0000-0003-3770-0513

627 KU Leuven

628 filip.buekens@kuleuven.be

## 629 References

- 630 ÁSTA [ÁSTA KRISTJANA SVEINSDÓTTIR]. 2018. *Categories We Live By: The Construction*  
 631 *of Sex, Gender, Race, and Other Social Categories*. Oxford: Oxford University Press,  
 632 doi:10.1093/oso/9780190256791.001.0001.
- 633 EUBANK, Nicholas. 2015. “Taxation, Political Accountability and Foreign Aid: Lessons  
 634 from Somaliland.” *The Journal of Development Studies* 48(4): 465–480, doi:10.108  
 635 0/00220388.2011.598510.
- 636 FLANIGAN, Shawn T. 2014. “Motivations and Implications of Community Service  
 637 Provision by La Familia Michoacána / Knights Templar and other Mexican Drug  
 638 Cartels.” *Journal of Strategic Security* 7(3): 63–83, doi:10.5038/1944-0472.7.3.4.

---

24 The two problems, i.e., constrained revisability and responsiveness to syntactic properties, are, of course, linked in that both are the result of the fact that our interests in making regulative judgments are practical and normative. I would like to thank an anonymous referee for highlighting this point.

\* THANKS



- 639 GRIER, David Alan. 2005. *When Computers Were Human*. Princeton, New Jersey:  
640 Princeton University Press, doi:10.1515/9781400849369.
- 641 GUALA, Francesco. 2015. "The Normativity of Institutions." *Phenomenology and Mind*  
642 9: 118–128, doi:10.13128/Phe\_Mi-18157.
- 643 —. 2016. *Understanding Institutions: The Science and Philosophy of Living Together*.  
644 Princeton, New Jersey: Princeton University Press, doi:10.1515/9781400880911.
- 645 GUALA, Francesco and HINDRIKS, Frank. 2015. "A Unified Social Ontology." *The*  
646 *Philosophical Quarterly* 65(259): 177–201, doi:10.1093/pq/pqu072.
- 647 HART, H. L. A. 1961. *The Concept of Law*. Oxford: Oxford University Press. Third  
648 edition: Hart (2012).
- 649 —. 2012. *The Concept of Law*. 3rd ed. Oxford: Oxford University Press. First edition:  
650 Hart (1961), third edition with a postscript, edited by Penelope A. Bulloch and  
651 Joseph Raz, doi:10.1093/he/9780199644704.001.0001.
- 652 HINDRIKS, Frank and GUALA, Francesco. 2015. "Institutions, Rules, and Equilibria: A  
653 Unified Theory." *Journal of Institutional Economics* 11(3): 459–480, doi:10.1017/S1  
654 744137414000496.
- 655 QUINE, Willard van Orman. 1951. "Two Dogmas of Empiricism." *The Philosophical*  
656 *Review* 60(1): 20–43. Reprinted in Quine (1953, 20–46), doi:10.2307/2181906.
- 657 —. 1953. *From a Logical Point of View: 9 Logico-Philosophical Essays*. Cambridge,  
658 Massachusetts: Harvard University Press. Cited after the revised edition: Quine  
659 (1961).
- 660 —. 1961. *From a Logical Point of View: 9 Logico-Philosophical Essays*. 2nd ed. Cam-  
661 bridge, Massachusetts: Harvard University Press. Revised edition of Quine (1953),  
662 reprinted 1980.
- 663 RANSEDELL, Joseph. 1971. "Constitutive Rules and Speech-Act Analysis." *The Journal*  
664 *of Philosophy* 68(13): 385–400, doi:10.2307/2025037.
- 665 SEARLE, John R. 1995. *The Construction of Social Reality*. London: The Free Press.
- 666 —. 2003. "Reply to Barry Smith (2003)." *American Journal of Economics and Sociology*  
667 62(1): 299–309, doi:10.1111/1536-7150.t01-i-00012.
- 668 —. 2005. "What Is an Institution?" *Journal of Institutional Economics* 1(1): 1–22,  
669 doi:10.1017/S1744137405000020.
- 670 —. 2010. *Making the Social World: The Structure of Human Civilization*. Oxford: Oxford  
671 University Press, doi:10.1093/acprof:osobl/9780195396171.001.0001.
- 672 SMIT, J. P., BUEKENS, Filip and DU PLESSIS, Stan. 2011. "What is Money? An Al-  
673 ternative to Searle's Institutional Facts." *Economics and Philosophy* 27(1): 1–22,  
674 doi:10.1017/s0266267110000441.
- 675 —. 2014. "Developing the Incentivized Action View of Institutional Reality." *Synthese*  
676 191(8): 1813–1830, doi:10.1007/s11229-013-0370-5.
- 677 —. 2016. "Cigarettes, Dollars and Bitcoins—an Essay on the Ontology of Money." *Journal of Institutional Economics* 12(2): 327–347, doi:10.1017/S174413741500040  
678 5.  
679

- 680 SMITH, Barry. 2003. "The Ontology of Social Reality." *American Journal of Economics*  
681 *and Sociology* 62(1): 283–299, doi:[10.1111/1536-7150.t01-1-00012](https://doi.org/10.1111/1536-7150.t01-1-00012).
- 682 TURNER, Stephen P. 2010. *Explaining the Normative*. Cambridge: Polity Press.
- 683 WILSON, Robert A. 2007. "Social Reality and Institutional Facts: Sociality Within and  
684 Without Intentionality." in *Intentional Acts and Institutional Facts: Essays on John*  
685 *Searle's Social Ontology*, edited by Savas L. TSOHATZIDIS, pp. 139–153. Dordrecht:  
686 Springer, doi:[10.1007/978-1-4020-6104-2\\_6](https://doi.org/10.1007/978-1-4020-6104-2_6).

PROOF

# The Minimalist Theory of Truth and the Generalisation Problem

LI ZHANG & LEON HORSTEN

The Minimalist Theory of Truth must show how it can prove certain truth-involving generalisations. Horwich has proposed two solutions to this challenge over the past decades. The first of these invokes Hilbert's  $\omega$ -rule and is unacceptable. The second proposal can naturally be viewed in different ways. We show how this second proposal is naturally interpreted as a suggestion to solve the truth generalisation problem using uniform reflection rules. We also argue that this is indeed the right way for Horwich to respond to the truth generalisation problem.

Over the past three decades, Horwich's *minimalism* has been the most discussed deflationary truth theory. Generally speaking, this theory claims that everything about truth can be explained by the collection of underived and unproblematic instances of the equivalence schema

$\langle \text{ES} \rangle \langle p \rangle$  is true iff  $p$ .

In the literature, the equivalence schema  $\langle \text{ES} \rangle$  is also known as the Tarski-schema or T-schema; its instances are known as Tarski-biconditionals or T-sentences. The theory consisting of all underived, unproblematic Tarski-biconditionals, namely, the theory taking all such biconditionals to be *axioms*, is called the "*Minimalist Theory of Truth*" (MT).

Firstly, Horwich believes that truth is *non-substantial*, so we should not define truth with any *substantial* concept. Instead, the meaning of 'is true' is given by the collection of underived, unproblematic instances of the T-schema. Horwich believes that

our understanding of 'is true'—our knowledge of its meaning—consists in the fact that the explanatorily basic regularity in our use of it is the inclination to accept instantiations of the schema

714 (E) ‘the proposition that  $p$  is true if and only if  $p$ ,

715 by declarative sentences of English (including any extensions of  
716 English). (Horwich 1998, 35, emphasis omitted)

717 Due to its non-substantiality, truth should remain neutral in debates in other  
718 philosophical and non-philosophical areas (Horwich 1998, 52).

719 Secondly, Horwich argues that MT alone suffices to explain all the truth-  
720 involving phenomena (Horwich 1998, 5). Thus, for instance, without equating  
721 truth with utility or any other substantial concept, MT suffices to explain that  
722 “true beliefs help us to achieve practical goals” (Horwich 1998, 44). In what  
723 follows, we denote the first point Horwich has made as *the neutrality thesis*  
724 and the second as the *adequacy thesis* of minimalism (Gupta 1993, 361).

725 Despite Horwich’s clever arguments for the two minimalistic theses, many  
726 logicians and philosophers insist that Horwich’s minimalism is far from cor-  
727 rect since there are many truth-involving facts that cannot be explained by it.  
728 In particular, it cannot prove truth generalisations that we regard as acceptable.  
729 For instance, it is unclear how MT proves

730 (1) Every proposition of the form  $p \rightarrow p$  is true,

731 or

732 (2) Every proposition is such that either it or its negation is true.

733 In fact, many believe that it is impossible for MT to prove sentences such as  
734 (1) and (2). In the literature, this problem is known as *the truth generalisation*  
735 *problem* (Halbach 2014, 57; Raatikainen 2005, 177). Horwich has formulated  
736 two proposals in response to this challenge, but, as they stand, neither of them  
737 decisively answers the problem. We will defend an amplification and extension  
738 of Horwich’s second proposal and argue that this successfully tackles the truth  
739 generalisation problem within the framework of truth-theoretic minimalism.

740 The structure of this paper is as follows: in section 1, we reformulate Hor-  
741 wich’s minimalist truth theory in such a way that some unclarities of his  
742 original formulation are removed. In section 2, we show why MT and its  
743 modifications cannot prove intuitively acceptable truth generalisations. In  
744 section 3, we evaluate Horwich’s two proposals in the light of critiques of  
745 them that have appeared in the literature. In section 4, we show how MT  
746 proves an ample collection of truth generalisations when strengthened with

747 uniform reflection rules, and we argue this to be in line with Horwich’s second  
748 proposal. In section 5, we conclude this paper by suggesting that Horwich  
749 should accept our formulation of the reflection rules proposal since it coheres  
750 best with his other truth-theoretic theses.

## 751 **1 Reformulating MT**

752 It has been recognized that several aspects of the formulation of MT are  
753 unclear. In particular, it is not clear which Tarski-biconditionals belong to  
754 MT’s axioms and which do not. Moreover, it is not clear how taking proposi-  
755 tions as primary truth-bearers increases MT’s proof-theoretic strength. Thus,  
756 we suggest two modifications of MT in this paper. First, by applying the T-  
757 schema to sentences *that themselves do not contain the truth predicate*, we  
758 obtain a precise description of MT’s axioms. Second, we take sentences to be  
759 primary truth-bearers. Given these two modifications, MT is equivalent to  
760 the axiomatic truth theory TB (for “Tarski-biconditionals”) when we take the  
761 Peano Arithmetic to be its base theory.<sup>1</sup>

762 One reason for MT’s vagueness is Horwich’s approach to truth-theoretic  
763 paradoxes. Horwich concedes that if some instances of the T-schema are  
764 included in MT’s axioms, MT proves contradictions. He demonstrates this in  
765 a familiar way by applying the T-schema to the sentence

766 THE PROPOSITION FORMULATED IN CAPITAL LETTERS IS  
767 NOT TRUE.

768 He argues that the only acceptable strategy for this problem is to exclude some  
769 instances of the T-schema from the axioms of MT (Horwich 1998, 40–41). The  
770 spirit of his approach to paradoxes has been shared by prominent logicians,  
771 including Tarski: by putting different constraints on the scope of the T-schema,  
772 we obtain different formal truth theories. These theories capture central uses  
773 of the truth predicate while, in the meantime, being adroit at avoiding truth-  
774 theoretic paradoxes. What renders Horwich’s strategy different is that he does  
775 not give a specification of either the permitted or prohibited instances of the  
776 T-schema; he only requires that the collection of MT’s axioms should be a  
777 maximally consistent set of sentences (Horwich 1998, 42). Unfortunately,  
778 McGee has shown there are uncountably many mutually incompatible sets

---

1 TB is also sometimes denoted as DT (for “disquotational theory”) in the literature (Halbach 2014, 53).

779 that satisfy this requirement; none of them are recursively axiomatisable.  
 780 Therefore, Horwich must impose more constraints on the instances of the  
 781 T-schema (McGee 1992, 236–237).

782 TB is axiomatisable and consists of unproblematic Tarski-biconditionals,  
 783 which renders it a suitable substitute for MT. However, far be it from us to  
 784 claim that TB is the *only* suitable substitute for MT. Many natural axiomatic  
 785 disquotational theories of truth would do just as well. For instance, if one  
 786 would substitute variants of Halbach’s theory of *Positive Tarski-biconditionals*  
 787 (2009) for MT instead, then the arguments of the present article would still  
 788 go through.

789 Most logicians who are interested in formal truth theories, such as Tarski  
 790 (1944, 342), McGee (1992, 235), Halbach (2014, 12), and Cieśliński (2018,  
 791 1083, n.8), take sentences to be primary truth-bearers. The reasons for their  
 792 choice are quite straightforward: propositions are ill-understood and contro-  
 793 versial, whereas we have rigorous and widely accepted syntactical theories of  
 794 sentences.

795 Horwich nonetheless insists on formulating minimalism in terms of propo-  
 796 sitions because he believes that there exist propositions that cannot be ex-  
 797 pressed by current human languages:

798 Patrick Grim pointed out to me that the minimal theory cannot be  
 799 regarded as *the set* of propositions of the form  $\langle \langle p \rangle \text{ is true iff } p \rangle$ ;  
 800 for there is no such set. The argument for this conclusion is that  
 801 if there were such a set, then there would be distinct propositions  
 802 regarding *each* of its subsets, and then there would have to be  
 803 distinct axioms of the theory corresponding to those propositions.  
 804 Therefore there would be a 1 – 1 function correlating the subsets  
 805 of MT with some of its members. But Cantor’s diagonal argument  
 806 shows that there can be no such function. Therefore, MT is not  
 807 a set. In light of this *result* [our emphasis], when we say things  
 808 like “ $\langle A \rangle$  follows from the minimal theory,” we must take that to  
 809 mean, not that the relation of *following from* holds between  $\langle A \rangle$   
 810 and a certain entity, the minimal theory; but rather that it holds  
 811 between  $\langle A \rangle$  and *some part* of the minimal theory—i.e., between  
 812  $\langle A \rangle$  and some set of propositions of the form  $\langle \langle p \rangle \text{ is true iff } p \rangle$ .  
 813 (Horwich 1998, 20–21, n.4)

814 The particular argument of Grim that is alluded to here goes as follows (see  
 815 Grim 1988): Suppose there were a set  $S$  of all truths, and consider all subsets

816 of  $S$ , i.e., all members of the power set  $\mathcal{P}(S)$ . To each element of this power  
817 set will correspond a truth. To each element of the power set, for example, a  
818 particular truth  $p$  either will or will not belong as a member. In either case,  
819 we will have a truth: that  $p$  is a member of that element or that it is not. There  
820 will then be at least as many truths as there are elements of the power set  
821  $\mathcal{P}(S)$ . But by Cantor's theorem, we know that the power set of any set will  
822 be larger than the original. There will then be more truths than there are  
823 members of  $S$ , and for any set of truths  $S$ , there will be some truth left out.  
824 There can therefore be no set of all truths.

825 The quotation by Horwich shows that he regards Grim's argument as defini-  
826 tive: he takes the conclusion of the argument as a philosophical *result*. But  
827 it is far from clear whether, in the absence of a detailed, widely accepted  
828 theory of propositions and their constituents, Grim's argument is persuasive.  
829 To give but one example of a worry that one might have here,<sup>2</sup> observe that  
830 Grim's succinct argument presupposes that for each subset  $B$  of  $S$ , there exists  
831 a *proposition* of the form  $p \in S (\neg p \in S)$ . For this to be the case, for each  
832 such subset  $B$ , there has to be an individual concept of  $B$  as a part of this  
833 proposition. But whether all such individual concepts exist is a substantial  
834 and unsettled philosophical question.

835 In view of this, there seems to be no pressing need for Horwich to take  
836 truth to be a property of propositions. Nonetheless, we do not ask of Horwich  
837 that he abandon his views of the kinds of entities that are the bearers of truth.  
838 A sub-class of the totality of all propositions is the *set* of all propositions that  
839 can be expressed by sentences belonging to some fixed language. The theory  
840 of true *sentences* (of some language) can then be seen as a special case of  
841 Horwich's more general theory of truth of propositions. So the argument  
842 that is developed in the subsequent sections intends to support the thesis  
843 that, *as far as truth of (propositions expressed by) sentences goes*, Horwich's  
844 arguments of the early 2000s concerning truth generalisations were at least a  
845 decade ahead of their time, albeit not fully fleshed out. That Horwich might  
846 well be sympathetic to such an interpretation of his views concerning truth  
847 generalisations is indicated by the passage in his *Truth* book, where he says  
848 that

849 ordinary language suggests that truth is a property of propositions,  
850 and that utterances, beliefs, assertions, etc., inherit their truth-

---

2 We do not have the space to go deeply into the literature that has been generated by Grim's argument.

851 like character from their relationship to propositions. However,  
 852 [previous considerations] show that this way of seeing things has  
 853 no particular explanatory merit. The truth-like conception for  
 854 each type of entity is equally minimalistic. And by assuming any  
 855 one of them we can easily derive the others. (Horwich 1998, 102)

## 852 The Truth Generalisation Problem

857 A non-trivial general claim of the form “Every  $x$  is  $\varphi$ ” cannot be proved by  
 858 a finite collection of premises, each of which asserts that  $a_i$  is  $\varphi$ , for some  $i$ ,  
 859 except if there is an additional premise that says that every object is one of  
 860 this finite number of  $a_i$ ’s. This also applies to MT in the sense that a general  
 861 claim of the form “For every sentence  $x$  of the form  $p \rightarrow p$ ,  $x$  is true,” for  
 862 example, cannot be proved in MT. Indeed, it has been *proved* that such a truth  
 863 generalisation cannot be proved in TB (= MT) (Halbach 2014, 56–57). Many  
 864 such truth generalisations appear to be conceptual truths about the concept  
 865 of truth. In particular, this is so for the classical compositional axioms of  
 866 truth that state that truth commutes with the logical connectives. Moreover,  
 867 there are valid philosophical and natural language arguments whose *validity*  
 868 depends not just on Tarski-biconditionals but also on compositional truth  
 869 axioms (Fujimoto 2022; Horsten 2010).

870 This poses a challenge to the minimalist theory of truth: recall that MT’s  
 871 adequacy thesis claims that *all* facts whose expression involves the truth  
 872 predicate can be explained by assuming no more about truth than instances  
 873 of the equivalence schema (Horwich 1998, 23). A number of philosophers  
 874 and logicians, including Armour-Garb (2010, 698), Gupta (1993, 363–364),  
 875 Halbach (2001, 1959–1960), and Soames (1997, 30–31), regard its inability to  
 876 prove truth generalisations as a serious defect of MT.

877 One may be tempted to appeal to “McGee’s trick” (McGee 1992, 238) and  
 878 contend that since it is always possible to find a T-sentence that is equivalent  
 879 to a given truth generalisation, when MT is not identified with TB but instead  
 880 with TB plus such additional Tarski-biconditionals, MT is capable of proving  
 881 all acceptable truth generalisations. Indeed, by the diagonal lemma, for every  
 882 truth generalisation  $A$ , there is a sentence  $\kappa$  such that

$$\vdash \kappa \leftrightarrow (T(\kappa) \leftrightarrow A).$$

883 By associativity of  $\leftrightarrow$ , the Tarski-equivalence  $\kappa \leftrightarrow T(\kappa)$  is provably equivalent  
 884 to  $A$ , where  $A$  is an acceptable truth generalisation. However, it is widely ac-



885 cepted that sentences such as that expressed by  $\kappa$  should not be allowed in  $\langle \text{ES} \rangle$   
886 since, by exactly the same procedure, it is possible to find a T-sentence equiv-  
887 alent to “Santa Claus exists,” which should not follow from any acceptable  
888 truth theory.

889 In sum, the generalisation problem poses a serious challenge for Horwich’s  
890 truth theory.

### 893 **3 Horwich’s Responses**

892 It is not clear exactly when Horwich came to the conclusion that MT cannot  
893 prove acceptable truth generalisations. But it is clear he wants to resolve this  
894 problem by strengthening MT with further theoretical resources. Moreover, it  
895 is possible to group his many responses into two categories: the  $\omega$ -rule proposal  
896 and a reflection-based proposal. We review these proposals in turn.

#### 371 *Horwich’s First Attempt: The $\omega$ -Rule*

898 In the postscript of the revised *Truth*, Horwich formulates his first attempt at  
899 solving the truth generalisation problem. There, he writes:

900 However, it seems to me that in the present case, where the topic is  
901 *propositions*, we can find a solution to this problem. For it is plau-  
902 sible to suppose that there is a truth-preserving rule of inference  
903 that will take us from a set of premises attributing to each proposi-  
904 tion some property,  $F$ , to the conclusion that all propositions have  
905  $F$ . No doubt this rule is not *logically* valid, for its reliability hinges  
906 not merely on the meanings of the logical constants, but also on  
907 the nature of propositions. But it is a principle we do find plau-  
908 sible. We commit ourselves to it, implicitly, in moving from the  
909 disposition to accept any proposition of the form “ $x$  is  $F$ ” (where  
910  $x$  is a proposition) to the conclusion “All propositions are  $F$ .” So  
911 we can suppose that this rule is what sustains the explanations  
912 of the generalizations about truth with which we are concerned.  
913 Thus we can, after all, defend the thesis that the basic theory of  
914 truth consists in some subset of the instances of the equivalence  
915 schema. (Horwich 1998, 137–138)

916 It has been acknowledged that the above-mentioned truth-preserving rule  
 917 amounts to a form of the  $\omega$ -rule (Raatikainen 2005, 175). Hilbert introduces  
 918 this principle in the following manner:

919 If it has been proved, for any given numeral  $\delta$ , that the formula

$$\mathfrak{A}(\delta)$$

920 is always a correct numerical formula, then the formula center

$$(x)\mathfrak{A}(x)$$

921 can be laid down as a starting formula [Ausgangsformel]. (Hilbert  
 922 1931, 1154)

923 Feferman rightly observed that Hilbert's own formulation of the  $\omega$ -rule is  
 924 somewhat vague (Feferman 1986, 212). The  $\omega$ -rule is perhaps more clearly  
 925 expressed as: "From infinitely many premises  $\varphi(0), \varphi(1), \dots$  that result from  
 926 replacing the numerical variable  $n$  in  $\varphi(n)$  with the numeral for each natural  
 927 number, conclude  $\forall x\varphi(x)$ " (Hazen 1998).

928 The  $\omega$ -rule is a strong rule: When enriched with this rule, PA proves true  
 929 arithmetic (Hazen 1998). With regard to the generalisation problem, when  
 930 augmented with the  $\omega$ -rule, MT is able to prove all acceptable truth general-  
 931 isations. Take a finite first-order language as an example; every sentence of  
 932 the form  $p \rightarrow p$  is a theorem of this language. Enumerate all sentences of the  
 933 form  $p \rightarrow p$ , so each of them is represented by a numeral. Apply T-sentences  
 934 of MT to them, so for each  $n$ ,  $T(n)$ . By the  $\omega$ -rule, we obtain the general claim  
 935  $\forall xT(x)$ .

936 However, certain features of the  $\omega$ -rule render this proposal problematic  
 937 and, in particular, unacceptable to the minimalist truth theory. Raatikainen  
 938 has argued that we, as finite human beings, cannot take infinitely many  
 939 premises into consideration simultaneously. Therefore, even if the theory  
 940 MT + the  $\omega$ -rule is capable of proving acceptable truth generalisations, those  
 941 generalisations are beyond the reach of ordinary human beings (Raatikainen  
 942 2005, 176). This problem with the  $\omega$ -rule cannot be overcome: it simply has  
 943 no effective (read: recursively enumerable) equivalent. Moreover, the proof-  
 944 theoretic strength of the  $\omega$ -rule makes it specifically unacceptable to the  
 945 minimalist truth theory. When enriched with this rule, Peano Arithmetic  
 946 proves all true arithmetic sentences. True arithmetic is not axiomatisable,  
 947 while MT is intended to be an axiomatised truth theory.

948 It is not clear whether or not Horwich has accepted critiques of his first  
 949 proposal. In a recent publication, Horwich still seems to propose using the  
 950  $\omega$ -rule as a solution to the truth generalisation problem:

951 For it is plausible to suppose that there is a truth-preserving rule  
 952 of inference that will take us from a set of premises attributing to  
 953 each proposition of a certain form some property,  $G$ , to the conclu-  
 954 sion that the *all* proposition have property  $G$ . And this rule—not  
 955 *logically* valid, but nonetheless necessarily truth-preserving given  
 956 the nature of proposition—enables the general facts about truth  
 957 to be explained by their instances. (Horwich 2003, 84, n.14, cited  
 958 after reprint)

959 Yet, in most of his recent writings, Horwich advocates an alternative resolution  
 960 based on an introspective process. To this proposal, we now turn.

### 3.2 Horwich's Second Attempt: Reflection

962 Over the years, Horwich's formulation of his second proposal has varied, and  
 963 it is not easy to select a preferred formulation from these variants. Extant  
 964 critiques of his various formulations are indecisive. Nonetheless, we will argue  
 965 that all variants of Horwich's second proposal need emendation in order to  
 966 solve the truth generalisation problem.

967 A first formulation of Horwich's second attempt emerges in (2001):

968 Whenever someone can establish, for any  $F$ , that it is  $G$ , and  
 969 recognizes that he can do this, then he will conclude that every  $F$   
 970 is  $G$ . (Horwich 2001, 157)

971 Call this *Solution 2.0*. This solution also consists in adding an additional rule  
 972 of inference to MT. But the additional rule of inference of *Solution 2.0* is  
 973 different from the  $\omega$ -rule.

974 In a revised version (2010) of the same paper, Horwich formulates a vari-  
 975 ant of this new proposal, which in effect amounts to a further, *substantially*  
 976 *different* proposal:

977 Whenever someone is disposed to accept, for any proposition of  
 978 structural type  $F$ , that it is  $G$  (and to do so for uniform reasons)  
 979 then he will be disposed to accept that every  $F$ -proposition is  $G$ .  
 980 (Horwich 2010, 45)

981 To the above statement, he adds the following *proviso*:

982 We cannot conceive of there being additional Fs—beyond those  
 983 Fs we are disposed to believe are G—which we would not have  
 984 the same sort of reason to believe are Gs. (Horwich 2010, 44)

985 Call the proposal that is encapsulated in the previous two quotations *Solution*  
 986 2.1. (Horwich endorses this same solution in 2003, 84.)

987 Armour-Garb argues that [Solution 2.1](#) is unsatisfactory because

988 one will not be disposed to accept (the proposition) that all F-  
 989 propositions are G, from the fact that, for any F-proposition, she  
 990 is disposed to accept that it is G (NB, even for uniform reasons),  
 991 unless she is *aware* of the fact that, for any F-proposition, she is  
 992 disposed to accept that it is G. (Armour-Garb 2010, 699)

993 The *proviso* Horwich added to [Solution 2.1](#) does not provide such an awareness  
 994 component. It merely adds a *negative* condition (“not being able to conceive of  
 995 there being F’s that are not G”), while Armour-Garb’s awareness requirement  
 996 is a positive condition. Nonetheless, [Solution 2.0](#) incorporates exactly the  
 997 awareness condition that Armour-Garb insists on (“and recognises that he  
 998 can do this”).

999 Armour-Garb is making a psychological observation here, but there is an  
 1000 accompanying *rational* point to be made also. If one does not *recognise* that  
 1001 for any *F*-proposition, she is disposed to accept that it is *G*, then she is not,  
 1002 without further ado, *rationally required* to believe that every *F*-proposition is  
 1003 *G*. *Ought* implies *can*, and, in this situation, she simply lacks the ground for  
 1004 accepting that every *F*-proposition is *G*.<sup>3</sup> For this reason, Horwich’s [Solution](#)  
 1005 [2.0](#) must be regarded as superior to his [Solution 2.1](#).

1006 Nonetheless, Armour-Garb would not be satisfied with [Solution 2.0](#) either.  
 1007 He argues that the switch, in the move from the premise to the conclusion of  
 1008 the rule of inference in [Solution 2.1](#), of “for any *F*-proposition” from outside  
 1009 the “disposed to accept”-context to inside the “disposed to accept”-context,  
 1010 is “viciously circular.” He is certainly right that this quantifier shift, which is  
 1011 also present in [Solution 2.0](#), is not derivable in classical logic. Nonetheless,  
 1012 we take issue with this aspect of Armour-Garb’s critique of Horwich’s second  
 1013 proposal. Indeed, we agree with Cieśliński that Armour-Garb’s dismissal of

---

3 Further discussion of these important matters can be found in Horsten (2021).

1014 Horwich’s second solution on the grounds of its being viciously circular is  
 1015 “hasty” (Cieśliński 2018, 1082): we will come back to this later.

1016 It is time to spell out the content of Horwich’s [Solution 2.0](#), i.e., the first  
 1017 quotation in this section, in more precise terms. We do this by formalising  
 1018 Horwich’s informally expressed—and somewhat vague—rule in first-order  
 1019 logic. In our formalisation of the first quotation in this section, we want to be  
 1020 charitable to Horwich. We do not claim that Horwich would agree with our  
 1021 formalisation (Horwich can speak for himself), but we will argue that there  
 1022 are good reasons for him to do so. Firstly, [Solution 2.0](#) contains the phrase  
 1023 “will conclude,” making it seem like a psychological prediction.<sup>4</sup> If it is taken  
 1024 in this way, then whether it is true or not is an empirical matter. But this is  
 1025 presumably not what Horwich intends. Rather, what he means is that the  
 1026 agent will be disposed to draw this conclusion *if she is rational*. In other words,  
 1027 Horwich purports to propose a rational *rule of inference* here. So it might  
 1028 be better to replace, in [Solution 2.0](#), “will conclude” with “may (rationally)  
 1029 conclude,” or perhaps even “should (rationally) conclude.” Secondly, since  
 1030 we are concerned with *establishing* truth generalisations, we identify the con-  
 1031 cepts “being disposed to accept” and “recognising” with being *provable*. In  
 1032 particular, we interpret the clause “and recognizes that he can do this” as *de re*  
 1033 provability of an arbitrary  $F$  that it is  $G$ . Thirdly, we identify provability with  
 1034 provability in the background theory, which is  $MT$ . If we were to identify prov-  
 1035 ability with provability in the system *including the rule*, then the proposed rule  
 1036 would indeed be viciously circular, confirming Armour-Garb’s (unfounded)  
 1037 suspicions. But if we identify provability with provability in  $MT$ , then there is  
 1038 no circularity. Fourthly, we *omit* the concept of provability (“being disposed  
 1039 to accept”) from the conclusion of the rule. With these precisifications in  
 1040 place—which we take to be reasonable, but we leave it open whether they are  
 1041 *exactly* in accordance with what Horwich intended—we obtain the following  
 1042 schematic rule:<sup>5</sup>

$$\frac{\vdash \forall x : F(x) \rightarrow Bew_{MT}(G(x))}{\vdash \forall x : F(x) \rightarrow G(x)}.$$

1043 We will call this rule H (for: “Horwich”). Observe that, unlike the  $\omega$ -rule, H is  
 1044 an *effective* rule: adding it to  $MT$  yields an axiomatic system.

1045 Worries based on the lottery paradox might cause one to doubt the ratio-  
 1046 nality of rule H. For any ticket (in a large, fair lottery), I believe that it is

4 Cieśliński sees this as the main weakness of Horwich’s recent views; see Cieśliński (2017, 80).

5 In the interest of readability, we are sloppy with Gödel coding here as well as later on in this article.

1047 not the winning ticket (and I believe this for “uniform reasons”). But from  
 1048 this, I am not prepared to infer that every ticket is a losing ticket (Kyburg  
 1049 1970, 56). Nonetheless, such a worry would be ill-founded, for the situation  
 1050 under consideration is different in one key respect. The irrationality of the  
 1051 lottery paradox inference stems from the fact that many small but non-zero  
 1052 probabilities (of being the winning ticket) can add up to a large probability  
 1053 (of one of a large collection of tickets being the winning one). But what is  
 1054 provable has probability 1 rather than  $1 - \epsilon$  (for some small  $\epsilon$ ) since provability  
 1055 in a sound system from necessary premises is itself necessary, and necessary  
 1056 truths by a Kolmogorov axiom for probability receive probability 1. So, the  
 1057 fair lottery phenomenon is irrelevant to the evaluation of rule H.<sup>6</sup>

#### 1054 **4 Uniform Reflection and Truth Generalisations**

1059 We have seen that Horwich recognises that H is not an admissible inference  
 1060 rule of first-order logic. The main questions that we want to answer in this  
 1061 section about H are the following: *To what extent and in which way does adding*  
 1062 *H to MT allow us to prove truth generalisations? Moreover: Is H a rational rule*  
 1063 *of inference?*

#### 4.1 *H and Uniform Reflection*

1065 It is clear that given a *sound* theory S, adding H (with  $Bew_{MT}$  replaced by  
 1066  $Bew_S$ ) to S results in a sound system. So, in particular,  $MT + H$  is a sound  
 1067 system.

1068 Next, we make the crucial observation that H is equivalent to a reflection  
 1069 rule that has been intensively investigated in proof theory. To this end, we  
 1070 first recall the notion of *uniform reflection principle* for a theory S (denoted as  
 1071  $RFN(S)$ ),

$$\forall x : Bew_S(\varphi(x)) \rightarrow \varphi(x),$$

1072 and the notion of *uniform reflection rule* for a theory S (denoted as  $UR_S$ ),

$$\frac{\vdash \forall x : Bew_S(\varphi(x))}{\vdash \forall x : \varphi(x)}.$$

---

6 An extended discussion of the relevance or irrelevance of the lottery paradox in this context can be found in Cieśliński (2017, sec. 13.5).

1073 Feferman has proved the remarkable little fact that  $\text{RFN}(S)$  is equivalent to  
 1074  $\text{UR}_S$  (Feferman 1962, Theorem 2.19). In light of this, it is easy to see that H  
 1075 is equivalent to  $\text{UR}_{\text{MT}}$  (and therefore also to  $\text{RFN}(\text{MT})$ ): the  $\Rightarrow$ -direction is  
 1076 obvious, and the  $\Leftarrow$ -direction follows immediately from Feferman’s theorem.

1077 At this point, a connection with Horwich’s *first* solution also becomes  
 1078 apparent. Indeed, the uniform reflection rule is widely seen as an effective  
 1079 version (a “tamed” version) of the  $\omega$ -rule. Horwich’s appeal to the  $\omega$ -rule was  
 1080 (rightly) rejected by Raatikainen on account of its non-effectiveness. Uniform  
 1081 reflection rules cannot be rejected on the same grounds.

1082 We will now see how the main observation of this subsection allows us to  
 1083 answer the question to what extent H enables us to prove truth generalisations.

#### 4.2 Deriving Truth Generalisations

1085 Let us denote  $\text{MT} + \text{H}$  as  $\text{MT}_1$ . Now that we have made Horwich’s [Solution](#)  
 1086 [2.0](#) precise, we address the question whether  $\text{MT}_1$  can prove all intuitively  
 1087 acceptable truth generalisations. An apparent counterexample is a proposition  
 1088 such as “There are as many truths as there are untruths” (Gupta 1993, 363).  
 1089 But this is a second-order statement involving not just sentences but also *sets*  
 1090 of sentences. So, it falls outside the scope of  $\text{MT}$  (=  $\text{TB}$ ), which cannot even  
 1091 express *claims* involving sets of sentences.

1092 The truth theory that takes the axioms that state that truth commutes  
 1093 with the logical connectives for sentences that do not themselves contain  
 1094 the notion of truth is called *CT*. It is fairly generally accepted that in *CT*, a  
 1095 vast amount of intuitively acceptable truth generalisations logically follow  
 1096 (Horsten 2011, chap. 6). So if Horwich can derive the truth axioms of *CT*, then  
 1097 he has made significant progress towards solving the truth generalisation  
 1098 problem. Nonetheless, it would be an exaggeration to say that *all* intuitively  
 1099 acceptable truth generalisations are provable in *CT*:<sup>7</sup> the truth generalisation  
 1100 “All arithmetical theorems of *CT* are true,” for instance, is not provable even  
 1101 in *CT*.

1102 With only one exception, the compositional truth axioms of *CT* can indeed  
 1103 be derived in  $\text{MT}_1$  (Horsten and Leigh 2017). As an example, let us consider  
 1104 the compositional axiom for negation:

$$\forall x \in \mathcal{L}_{\text{PA}} : T(\neg x) \leftrightarrow \neg Tx.$$

---

7 Thanks to an anonymous referee for making this point.

1105 Every *instance* of this axiom can be proved in TB (using Tarski-biconditionals).  
 1106 Moreover, *that* every instance can be proved in TB can be uniformly recognised  
 1107 (i.e., proved) as a combinatorial fact even in the background theory PA. So we  
 1108 have:

$$PA \vdash \forall x \in \mathcal{L}_{PA} : Bew_{MT}(T(\neg x) \leftrightarrow \neg Tx).$$

1109 Then, by UR<sub>MT</sub>, we indeed obtain  $\forall x \in \mathcal{L}_{PA} : T(\neg x) \leftrightarrow \neg Tx$ .

1110 The other compositional axioms can be derived in a similar way in MT<sub>1</sub>,  
 1111 with the sole exception of the quantifier axiom:

$$\forall \varphi(x) \in \mathcal{L}_{PA} : T(\forall x \varphi(x)) \leftrightarrow \forall x T\varphi(x).$$

1112 We cannot prove in MT, for every  $\varphi(x) \in \mathcal{L}_{PA}$ , that  $T(\forall x \varphi(x)) \leftrightarrow \forall x T\varphi(x)$ .  
 1113 The reason is that TB (= MT) only contains Tarski-biconditionals for *sen-*  
 1114 *tences*, i.e., for *closed* formulas. In order to prove, for each  $\varphi(x) \in \mathcal{L}_{PA}$ ,  
 1115 that  $T(\forall x \varphi(x)) \leftrightarrow \forall x T\varphi(x)$ , we need a slight strengthening of the Tarski-  
 1116 biconditionals of TB, namely, the *uniform* arithmetical Tarski-biconditionals,  
 1117 which are the sentences of the form  $\forall x (T\varphi(x) \leftrightarrow \varphi(x))$  for formulas  $\varphi(x) \in$   
 1118  $\mathcal{L}_{PA}$ . The resulting slight strengthening of TB is called *UTB*.

1119 How do we derive these uniform Tarski-biconditionals? We can prove them  
 1120 in MT<sub>1</sub> as follows (Horsten and Leigh 2017, Theorem 9).<sup>8</sup> Every instance of a  
 1121 given uniform (arithmetical) Tarski-biconditional can be proved in TB. This  
 1122 combinatorial fact can again be proved even in PA:

$$PA \vdash \forall x \in \mathcal{L}_{PA} : Bew_{MT} : T\varphi(x) \leftrightarrow \varphi(x).$$

1123 So, by applying UR<sub>MT</sub> in MT<sub>1</sub> to this fact, we obtain the result. Now, in a  
 1124 second stage, we can proceed as we did with the negation axiom. But to carry  
 1125 out this proof, we need to appeal to UR<sub>MT<sub>1</sub></sub>, which is the uniform reflection  
 1126 rule for MT<sub>1</sub>:

$$\frac{\vdash \forall x : Bew_{MT_1}(\varphi(x))}{\vdash \forall x : \varphi(x)},$$

1127 where  $\varphi$  can be any arithmetical formula, and  $Bew_{MT_1}$  formally expresses  
 1128 provability in MT<sub>1</sub>. For the same reasons as why UR<sub>MT</sub> exceeds MT, the rule

---

8 Theorem 9, obtained in Horsten and Leigh (2017), is based on uniform reflection principles rather than rules, but we have seen above that by an argument due to Feferman, the two are provably equivalent.



1129 UR<sub>MT<sub>1</sub></sub> exceeds MT<sub>1</sub>. If we apply this inference rule to the earlier obtained  
 1130 fact that PA proves

$$\forall x \in \mathcal{L}_{PA} : Bew_{MT_1}(T(\forall x\varphi(x)) \leftrightarrow \forall xT\varphi(x)),$$

1131 then we obtain the desired result that  $\forall x \in \mathcal{L}_{PA} : T(\forall x\varphi(x)) \leftrightarrow \forall xT\varphi(x)$ .

1132 In sum, we can prove all the compositional truth axioms of CT and, there-  
 1133 fore, many intuitively acceptable truth generalisations in  $MT_2 = MT_1 +$   
 1134  $UR_{MT_1} = MT + UR_{MT} + UR_{MT+UR_{MT}}$  (Horsten and Leigh 2017). In other  
 1135 words, many truth generalisations follow from two iterations of uniform re-  
 1136 flection on MT. Even more truth generalisations can be proved when this  
 1137 strategy is iterated further. By adding further uniform reflection principles to  
 1138 MT<sub>2</sub>, for instance, also the truth generalisation “All arithmetical theorems of  
 1139 CT are true” becomes provable.

1140 At this point, we see that we have to go slightly beyond our charitable  
 1141 interpretation of Horwich’s Solution 2.0. Horwich claims that *one* level of  
 1142 reflection on MT suffices to prove all acceptable truth generalisations. We  
 1143 now see that *two* levels of reflection on MT are required. Given the equiva-  
 1144 lence between Horwich’s rule H and Feferman’s uniform reflection rule, all  
 1145 acceptable truth generalisations can be derived in the theory  $MT + H + H'$ ,  
 1146 where H’ is just like H, except that its background theory is  $MT + H$  instead  
 1147 of MT:

$$\frac{\vdash \forall x : F(x) \rightarrow Bew_{MT+H}(G(x))}{\vdash \forall x : F(x) \rightarrow G(x)}.$$

1148 In sum, if H and H’ are *rational* rules of inference, then Horwich was very  
 1149 much on the right track.

### 4.3 Rationality

1151 Uniform reflection rules are rules that contain the required “awareness” com-  
 1152 ponent in the antecedent (the agent has to have a proof) and that are also, *pace*  
 1153 Armour-Garb, not circular in any way. In addition, in the premise of uniform  
 1154 reflection rules, the awareness/recognition component that is required is *proof*  
 1155 from the Tarski-biconditionals.

1156 On our interpretation and emendation of his view, Horwich contends that it  
 1157 is *rational* to add UR<sub>MT</sub> and UR<sub>MT<sub>1</sub></sub> to MT. With this, he would not be alone. In  
 1158 his work on *implicit commitment*, Feferman claimed that if an agent explicitly  
 1159 accepts a theory S, then she also ought to accept uniform reflection principles

1160 and rules for S, such as  $UR_S$  and  $UR_{S+UR_S}$  (Feferman 1991, 2, 44). Acceptance  
 1161 of  $UR_S$  is, in his view, *implicit* in the acceptance of S, and acceptance of  
 1162  $UR_{S+UR_S}$  is “implicitly implicit” in the acceptance of S.

1163 Feferman did not give an epistemological argument for *why*, if one accepts  
 1164 a theory S, one should also accept  $UR_S$  (and  $UR_{S+UR_S}$ ). A recent attempt to  
 1165 provide such an argument is given by Fischer in (2023), which can, in retrospect,  
 1166 be seen as one attempt to develop Horwich’s [Solution 2.0](#) in detail. A discus-  
 1167 sion of Fischer’s argument is outside the scope of this article. Here, we restrict  
 1168 ourselves to a few remarks on the issue. The uniform reflection rule for the  
 1169 theory that one is currently working in expresses a form of trust or confidence  
 1170 in this theory. If the theory one is working in is justified, then this trust is  
 1171 also justified, and therefore, accepting the uniform reflection rule is justified.  
 1172 The theory that is relevant in the present context is the truth theory MT. Hor-  
 1173 wich argues that this theory is indeed justified because Tarski-biconditionals  
 1174 express the content or meaning of the concept of truth (Horwich 2010, 17).  
 1175 Therefore, making one’s trust in MT explicit by accepting  $UR_{MT}$  and  $UR_{MT_1}$   
 1176 is rational.<sup>9</sup> Since, by Feferman’s theorem, H is equivalent to  $UR_{MT}$ , and H’ is  
 1177 equivalent to  $UR_{MT_1}$ , H and H’ are, therefore, also rational inference rules.

## 1175 5 Horwich Vindicated?

1179 There have been two phases in the history of truth-theoretic deflationism. In  
 1180 the first phase, disquotational axioms were taken to express the full content  
 1181 of the concept of truth. This phase comprises, a.o., Quine’s views on truth as  
 1182 a tool for semantic ascent and descent (Quine 1970, 10–13) and the prosentential  
 1183 theory of truth (Grover, Camp and Belnap 1975). Horwich’s minimalism  
 1184 is often viewed as a late and particularly bright exponent of this phase of  
 1185 deflationism. In the second phase, compositional axioms were taken to ex-  
 1186 press basic properties of the concept of truth. This phase started sometime in  
 1187 the 1980s, partly under the influence of Davidson’s truth-conditional com-  
 1188 positional approach to natural language semantics (Davidson 1967). During  
 1189 much of this second phase, Horwich’s views on the concept of truth came to  
 1190 be increasingly seen as dated and untenable. As a result of this, his writings  
 1191 about the generalisation problem after the first edition of his book *Truth* were  
 1192 mostly ignored by the truth-theoretic community.

---

9 Considerations such as these may provide at least the beginnings of a response to Cieślński’s complaint above (*confer supra*, page XX) that Horwich’s theory is too psychological.

1193 Perhaps we now experience the dawn of a third phase in the history of truth-  
1194 theoretic deflationism, in which the relation between the concept of truth on  
1195 the one hand and reflection principles on the other hand plays a major role. In  
1196 particular, it is currently a hotly debated question whether by making use of  
1197 reflection principles or rules, disquotationalism can solve the generalisation  
1198 problem. We make no attempt to adjudicate this discussion here. But we have  
1199 seen that Horwich anticipated the current philosophical debate already in  
1200 the early 2000s. So rather than being a truth-theoretic dinosaur, at the time  
1201 Horwich's views were ahead of their time—which, of course, does not mean  
1202 that they are in any way definitive.

1203 The main reason why Horwich's thoughts about the relation between re-  
1204 flection principles and truth generalisations were ignored is that Horwich's  
1205 view about this problem was not completely precise and was connected to  
1206 other views of his that can be separated from the problem at issue. Horwich  
1207 was committed to propositions as the bearers of truth but did not give a pre-  
1208 cise theory of propositions. At the same time, he was also committed to the  
1209 background disquotational theory as a maximal consistent collection of propo-  
1210 sitions, which prevents it from being recursively axiomatisable and, therefore,  
1211 prevents it from being learnable. But we have seen that a *derived* notion of  
1212 true proposition *expressible in a given language* makes perfect sense in Hor-  
1213 wich's framework. Moreover, Horwich's requirement of MT being a maximal  
1214 consistent collection of propositions is unrelated to his solution proposal to  
1215 the generalisation problem and can, therefore, simply be rejected—which is  
1216 exactly what the truth-theoretic community has largely done. In sum, Hor-  
1217 wich's views from the early 2000s on the truth generalisation problem can be  
1218 disentangled from the further commitments and unclarities with which he  
1219 connected them.

1220 The imprecision of his treatment of the generalisation problem prevented  
1221 Horwich from working out the technical details with full precision. For in-  
1222 stance, he did not see that *two* rounds of uniform reflection are needed in order  
1223 to derive the compositional truth principles from the disquotational axioms.  
1224 Nonetheless, Horwich did see that his strategy for dealing with the general-  
1225 isation problem is in line with his two main minimalistic theses: the neutrality  
1226 thesis and the adequacy thesis. Reflection rules are not truth-theoretic (or  
1227 *philosophical*) but *mathematical* rules. Uniform reflection rules are universally  
1228 seen as mathematical rules because they have substantial mathematical con-  
1229 sequences: they are canonical ways of extending the mathematical strength  
1230 of a theory. Therefore, strengthening MT with uniform reflection rules does

not affect the neutrality of the theory of *truth*. (Indeed, as mentioned earlier, MT can be taken to be proof-theoretically conservative over its background theory PA.) Moreover, since CT is derivable from MT by means of two rounds of uniform reflection, and CT proves the needed truth generalisations, a solution to the generalisation problem is reached, whereby the challenge to the adequacy thesis is answered.

The more recent debate about the connection between reflection principles and the truth generalisation problem developed only after 2015, and it developed largely independently from Horwich's views on the generalisation problem. Moreover, we now see further and more clearly in these matters than Horwich did around 2002. Yet it would be a mistake to take Horwich's early thoughts on this issue to be merely of historical relevance ("give credit where credit is due"). The appeal to proof-theoretic reflection principles and rules as a means to derive compositional truth axioms is sometimes seen as a mere "technical" manoeuvre. But Horwich, at the time, did not know any of the proof-theoretic literature concerning reflection principles and hit on the basic idea *in tempore non suspecto*. Purely by philosophically thinking about how to solve the generalisation problem in a disquotational framework, he, in one of his proposals (Solution 2.0), arrived at the view that the compositionality or truth follows by the uniform reflection rule from disquotational principles. This is simply amazing, and it shows that rather than being merely a technical trick, it is a very natural theoretical view to take.\*

Li Zhang

 0000-0002-7766-7263

Tsinghua University

l-zhang17@mails.tsinghua.edu.cn

Leon Horsten

 0000-0003-3610-9318

University of Konstanz

Leon.Horsten@uni-konstanz.de

---

\* The authors are indebted to the members of the Antos-Horsten doctoral seminar at the University of Konstanz, where an early version of this article was presented, for helpful and thoughtful comments. They also thank two anonymous referees, whose comments have led to substantial improvements of this article. The first author's research for this article was supported by a Chinese Scholarship Council Fellowship (CSC No. 201906210194) for carrying out doctoral research at the University of Bristol and at the University of Konstanz. Without the support of this scholarship, the present article could not have been written.

## References

- 1261
- 1262 ARMOUR-GARB, Bradley. 2010. "Horwichian Minimalism and the Generalization  
1263 Problem." *Analysis* 70(4): 693–703, doi:10.1093/analys/anq073.
- 1264 BEALL, J. C. and ARMOUR-GARB, Bradley, eds. 2005. *Deflationism and Paradox*. Oxford:  
1265 Oxford University Press, doi:10.1093/oso/9780199287116.001.0001.
- 1266 CIEŚLIŃSKI, Cezary. 2017. *The Epistemic Lightness of Truth: Deflationism and Its Logic*.  
1267 Cambridge: Cambridge University Press, doi:10.1017/9781108178600.
- 1268 —. 2018. "Minimalism and the Generalisation Problem: On Horwich's Second Solu-  
1269 tion." *Synthese* 195(3): 1077–1101, doi:10.1007/s11229-016-1227-5.
- 1270 DAVIDSON, Donald. 1967. "Truth and Meaning." *Synthese* 17(1): 304–323. Reprinted  
1271 in Davidson (1984, 17–36), doi:10.1007/BF00485035.
- 1272 —. 1984. *Inquiries into Truth and Interpretation*. Oxford: Oxford University Press,  
1273 doi:10.1093/0199246297.001.0001.
- 1274 EWALD, William Bragg, ed. 1996. *From Kant to Hilbert: A Source Book in the Foundations*  
1275 *of Mathematics. Volume II*. Oxford: Oxford University Press, doi:10.1093/oso/9780
- 1276 198505365.001.0001.
- 1277 FEFERMAN, Solomon. 1962. "Transfinite Recursive Progressions of Axiomatic Theo-  
1278 ries." *The Journal of Symbolic Logic* 27(3): 259–316, doi:10.2307/2964649.
- 1279 —. 1986. "Introductory Note to 1931c [Gödel (1931)]." in *Collected Works, Volume I:*  
1280 *Publications 1929–1936*, pp. 208–212. Oxford: Oxford University Press. Edited by
- 1281 Solomon Feferman, John W. Dawson, Jr., Stephen C. Kleene, Gregory H. Moore,  
1282 Robert M. Solovay and Jean van Heijenoort, doi:10.1093/oso/9780195147209.003.
- 1283 0019.
- 1284 —. 1991. "Reflecting on Incompleteness." *The Journal of Symbolic Logic* 56(1): 1–49,  
1285 doi:10.2307/2274902.
- 1286 FISCHER, Martin H. 2023. "Another Look at Reflection." *Erkenntnis* 88(2): 479–509,  
1287 doi:10.1007/s10670-020-00363-9.
- 1288 FUJIMOTO, Kentaro. 2022. "The Function of Truth and the Conservativeness Argu-  
1289 ment." *Mind* 131(521): 129–157, doi:10.1093/mind/fzaa083.
- 1290 GÖDEL, Kurt. 1931. "Besprechung von Hilbert (1931)." *Zentralblatt für Mathematik*  
1291 *und ihre Grenzgebiete* 1: 260. English translation in van Heijenoort (1967, 596–616);  
1292 reprinted in Gödel (1986, 213/214).
- 1293 —. 1986. *Collected Works, Volume I: Publications 1929–1936*. Oxford: Oxford University  
1294 Press. Edited by Solomon Feferman, John W. Dawson, Jr., Stephen C. Kleene,  
1295 Gregory H. Moore, Robert M. Solovay and Jean van Heijenoort, doi:10.1093/oso/
- 1296 9780195147209.001.0001.
- 1297 GRIM, Patrick. 1988. "Logic and Limits of Knowledge and Truth." *Noûs* 22(3): 341–367,  
1298 doi:10.2307/2215708.
- 1299 GROVER, Dorothy L. 1992. *A Prosentential Theory of Truth*. Princeton, New Jersey:  
1300 Princeton University Press, doi:10.1515/9781400862689.

- 1301 GROVER, Dorothy L., CAMP, Joseph L., Jr. and BELNAP, Nuel D., Jr. 1975. "A Prosen-  
1302 tential Theory of Truth." *Philosophical Studies* 27(2): 73–125. Reprinted as Grover  
1303 (1992, ch. 3), doi:10.1007/bf01209340.
- 1304 GUPTA, Anil. 1993. "Minimalism." in *Philosophical Perspectives 7: Language and Logic*,  
1305 edited by James E. TOMBERLIN, pp. 359–369. Atascadero, California: Ridgeview  
1306 Publishing Co., doi:10.2307/2214129.
- 1307 HALBACH, Volker. 2001. "Disquotational Truth and Analyticity." *The Journal of Sym-*  
1308 *bolic Logic* 66(4): 1959–1973, doi:10.2307/2694987.
- 1309 —. 2009. "Reducing Compositional to Disquotational Truth." *The Review of Symbolic*  
1310 *Logic* 2(4): 786–798, doi:10.1017/s1755020309990220.
- 1311 —. 2011. *Axiomatic Theories of Truth*. 1st ed. Cambridge: Cambridge University Press,  
1312 doi:10.1017/cbo9780511921049.
- 1313 —. 2014. *Axiomatic Theories of Truth*. 2nd ed. Cambridge: Cambridge University Press.  
1314 First edition: Halbach (2011), doi:10.1017/cbo9781139696586.
- 1315 HAZEN, Allen Patterson. 1998. "Non-Constructive Rules of Inference." in *The Routledge*  
1316 *Encyclopedia of Philosophy*, edited by Edward J. CRAIG. London: Routledge. The  
1317 Routledge Encyclopedia was made available online in 2002 and is now regularly  
1318 updated, doi:10.4324/9780415249126-Y014-1.
- 1319 VAN HEIJENOORT, Jan, ed. 1967. *From Frege to Gödel: A Source Book in Mathematical*  
1320 *Logic 1879–1931*. Cambridge, Massachusetts: Harvard University Press.
- 1321 HILBERT, David. 1931. "Die Grundlegung der elementaren Zahlenlehre." *Mathema-*  
1322 *tische Annalen* 104(1): 485–494. Translated as "The Grounding of Elementary  
1323 Number Theory" in Ewald (1996, 1148–1156), doi:10.1007/BF01457953.
- 1324 HORSTEN, Leon. 2010. "On a Necessary Use of Truth in Epistemology." in *The Ana-*  
1325 *lytical Way: Proceedings of the 6th European Congress of Analytic Philosophy*,  
1326 edited by Tadeusz CZARNECKI, Katarzyna KIJANIA-PLACEK, Olga POLLER, and  
1327 Jan WOLEŃSKI, pp. 371–376. Studies in Logic n. 25. London: College Publications.
- 1328 —. 2011. *The Tarskian Turn: Deflationism and Axiomatic Truth*. Cambridge, Mas-  
1329 sachusetts: The MIT Press, doi:10.7551/mitpress/9780262015868.001.0001.
- 1330 —. 2021. "On Reflection." *The Philosophical Quarterly* 71(4): 738–757, doi:10.1093/pq  
1331 /pqaa083.
- 1332 HORSTEN, Leon and LEIGH, Graham E. 2017. "Truth is Simple." *Mind* 126(501):  
1333 195–232, doi:10.1093/mind/fzv184.
- 1334 HORWICH, Paul. 1990. *Truth*. Oxford: Blackwell Publishers. Second edition: Horwich  
1335 (1998).
- 1336 —. 1998. *Truth*. 2nd ed. Oxford: Oxford University Press. First edition: Horwich (1990),  
1337 doi:10.1093/0198752237.001.0001.
- 1338 —. 2001. "A Defense of Minimalism." *Synthese* 126(1/2): 149–165. Substantially revised  
1339 version in Horwich (2010, 35–56), doi:10.1023/a:1005279406402.
- 1340 —. 2003. "A Minimalist Critique of Tarski on Truth." in *Philosophy and Logic in*  
1341 *Search of the Polish Tradition: Essays in Honour of Jan Woleński on the Occasion of*

- 1342 *His 60th Birthday*, edited by Jaakko HINTIKKA, Tadeusz CZARNECKI, Katarzyna  
1343 KIJANIA-PLACEK, Tomasz PLACEK, and Artur ROJSZCZAK, pp. 3–12. Synthese  
1344 Library n. 323. Dordrecht: Kluwer Academic Publishers. Reprinted, in substantially  
1345 revised form, in Beall and Armour-Garb (2005, 75–84) and in Horwich (2010,  
1346 79–97), doi:10.1007/978-94-017-0249-2\_1.
- 1347 —. 2010. *Truth – Meaning – Reality*. Oxford: Oxford University Press, doi:10.1093/acpr  
1348 of:oso/9780199268900.001.0001.
- 1349 KYBURG, Henry E., Jr. 1970. “Conjunctivitis.” in *Induction, Acceptance, and Rational*  
1350 *Belief*, edited by Marshall SWAIN, pp. 55–82. Synthese Library n. 26. Dordrecht:  
1351 D. Reidel Publishing Co., doi:10.1007/978-94-010-3390-9\_4.
- 1352 MCGEE, Vann. 1992. “Maximal Consistent Sets of Instances of Tarski’s Schema (T).”  
1353 *Journal of Philosophical Logic* 21(3): 235–241, doi:10.1007/bf00260929.
- 1354 QUINE, Willard van Orman. 1970. *Philosophy of Logic*. Foundations of Philosophy  
1355 Series. Englewood Cliffs, New Jersey: Prentice-Hall, Inc. Second edition: Quine  
1356 (1986).
- 1357 —. 1986. *Philosophy of Logic*. 2nd ed. Cambridge, Massachusetts: Harvard University  
1358 Press. First edition: Quine (1970).
- 1359 RAATIKAINEN, Panu. 2005. “On Horwich’s Way Out.” *Analysis* 65(3): 175–177, doi:10  
1360 .1111/j.1467-8284.2005.00546.x.
- 1361 SOAMES, Scott. 1997. “The Truth about Deflationism.” in *Philosophical Issues 8: Truth*,  
1362 edited by Enrique VILLANUEVA, pp. 1–44. Atascadero, California: Ridgeview Pub-  
1363 lishing Co., doi:10.2307/1522992.
- 1364 TARSKI, Alfred. 1944. “The Semantic Conception of Truth and the Foundations of  
1365 Semantics.” *Philosophy and Phenomenological Research* 4(3): 341–375, doi:10.230  
1366 7/2102968.

PROOF



# The Problem of Thomistic Parts

FR. JAMES DOMINIC ROONEY, OP

Thomas Aquinas embraces a controversial claim about the way in which parts of a substance depend on the substance's substantial form. The substantial form is responsible for the identity/nature of the parts of the substance such a form constitutes. Aquinas's controversial claim can be roughly put as the view that things are members of their kind in virtue of their substantial form.

The aim of this paper will be to defend Aquinas's claim that every time the  $x$ s come to compose a  $y$ , those  $x$ s have to undergo a change in kind membership. After defending the Thomistic account, I propose that approaching problems of material composition as a Thomist has a significant, oft-overlooked advantage of involving a thoroughgoing naturalistic methodology that resolves such problems by appeal to empirical considerations.

Thomas Aquinas embraces a controversial claim about the way in which parts of a substance depend on the substance's substantial form. On his metaphysics, a "substantial form" is not merely a relation among already existing things, in virtue of which (for example) the arrangement or configuration of those things would count as a substance. The substantial form is rather responsible for the identity or nature of the parts of the substance such a form constitutes (Marmodoro and Page 2016, 17–18). Substantial forms thus do not have substance-parts as that which they characterize, i.e., their matter. However, the implication is that if some substances come to compose another substance as proper parts, those things that become parts must ipso facto cease to be substances. Conversely, if a material part ceases to compose a substance as a part, that thing will become a substance or a heap of substances. Aquinas's controversial claim can be roughly put as the view that things are members of their kind in virtue of their substantial form. When a part ceases to compose a substance, it ceases to be that kind of thing that it was when it composed its parent substance and, so, loses all of the properties or powers that are associated with being a part of that substance.

1398 As an illustration of the implications of this claim, consider the death of  
 1399 Socrates. Aquinas holds that “the soul [...] is the form of the whole body and  
 1400 each of its parts. [...] Thus, it is necessary that each part of a man and that of  
 1401 an animal receive its existence and species from the soul as from its proper  
 1402 form”<sup>1</sup> (QDA,<sup>2</sup> a. 10 resp.). When Socrates dies, all of his parts—his body and  
 1403 his eyes and his skin—cease to have their act of existence and cease to be  
 1404 the things that they were when they composed Socrates. Socrates’s corpse  
 1405 does not have eyes or hands or skin, at least properly speaking, because, after  
 1406 the soul leaves the body, “neither eye nor flesh nor any part remains except  
 1407 equivocally”<sup>3</sup> (QDA, a. 10 resp.). To put it simply, Aquinas’s claim results in  
 1408 the implication that every time the *xs* come to compose a *y*, those *xs* have to  
 1409 undergo a change in kind membership (Koslicki 2008, 147).

1410 This has been called the “homonymy principle,” and it follows from  
 1411 Aquinas’s view of substantial forms and, specifically, from the position  
 1412 that substantial forms inform prime matter rather than substance-parts.  
 1413 Consequently, a substantial form must account for the determinate actuality  
 1414 of every part of the substance. Yet the homonymy principle has appeared  
 1415 to many to be so counterintuitive as to practically require a belief in the  
 1416 existence of substance-parts of substances. Kathrin Koslicki argues that if the  
 1417 homonymy principle were true, it would be impossible to explain continuity  
 1418 in change.<sup>4</sup> The aim of this paper will be to defend that the Thomistic claim  
 1419 that substantial forms account for the determinate actuality of every part of a  
 1420 substance is plausible and coherent. After defending the Thomistic account, I  
 1421 propose that approaching problems of material composition as a Thomist  
 1422 has a significant, oft-overlooked advantage of involving a thoroughgoing  
 1423 naturalistic methodology that resolves such problems by appeal to empirical  
 1424 considerations.

## 1425 **1 Being a Part of a Thomistic Substance**

1426 The Thomistic claim about substantial forms would not be controversial  
 1427 merely if it held that parts ceased to be parts when they ceased to compose a  
 1428 substance or that something became a part when it composed something else.

1 “anima [...] [e]st enim forma et totius corporis, et cuiuslibet partis eius. [...] Unde oportet quod quaelibet pars hominis et animalis recipiat esse et speciem ab anima sicut a propria forma.”

2 Thomas Aquinas, *Quaestiones Disputatae de Anima*.

3 “neque oculus neque caro neque aliqua pars remanet nisi aequivoce.”

4 This remains true in Koslicki (2018, 217–220).

1429 One could think, for this reason, there was a problem if we were to characterize  
1430 “substances” merely as those material objects that are not parts of any other.  
1431 What seems to be missing from my characterization of a substance is the way  
1432 in which a substance is a properly unified thing, as would result from having  
1433 a substantial form that actualizes all of that substance’s parts, making it one  
1434 *kind* of thing. Stump points out, for example, that on Aquinas’s view of what  
1435 it is to be a substance, “the ability to exist on its own is a necessary but not a  
1436 sufficient condition for something’s being a substance” (Stump 2003, 42).

1437 Stump attempts to appeal to the contemporary metaphysical concept of  
1438 emergence as that which sets apart Aquinas’s view of substance. But she  
1439 contrasts her view of emergence with other contemporary views because “on  
1440 Aquinas’s way of thinking about material objects what can emerge when form  
1441 is imposed on matter is not just properties but substances” (Stump 2003, 196).  
1442 She defines what it is to be an emergent whole, i.e., a substance: “W is an  
1443 emergent thing if and only if the properties and causal powers of W are not  
1444 simply the sum of the properties and causal powers of the constituents of W  
1445 when those constituents are taken singillatim, outside the configuration of  
1446 W” (Stump 2003, 43). Non-substances, like artifacts, are nothing more than  
1447 “the sum of [their] parts” because the properties and powers of such wholes  
1448 are nothing over and above the properties and powers of their parts (Stump  
1449 2003, 44).

1450 Obviously, Stump’s claim can be misleading without further qualification.  
1451 As Marmodoro and Page point out, emergence is overly permissive as a cri-  
1452 terion of substancehood: “There are plenty of examples of material objects  
1453 having—on account of their structure or external relations—emergent prop-  
1454 erties or functions that the parts individually do not have, without such objects  
1455 *ipso facto* being substances” (Marmodoro and Page 2016, 4). Yet Stump’s claim  
1456 is not that it is sufficient for something to be a substance if the parts actually  
1457 composing some substance lack properties or powers individually that the  
1458 whole substance possesses. Rather, Stump’s claim is that a substance has pow-  
1459 ers and properties that are not a sum of the powers and properties of the parts  
1460 that could *potentially* come to compose it. That is the sense of the qualification  
1461 that these parts must be considered apart from the actual configuration of the  
1462 whole.

1463 Similarly, Stump’s definition would be an insufficient characterization of  
1464 Aquinas’s views if it was understood as presuming this claim: that one and  
1465 the same thing can come to have properties or powers in virtue of composing  
1466 another substance as a proper part; or, that one and the same thing can lose

1467 powers or properties when ceasing to compose a whole and becoming a sub-  
 1468 stance. Such a view would violate Aquinas’s homonymy principle. Aquinas’s  
 1469 case of the death of Socrates showed that his body could not be identical  
 1470 when it was actually alive and when it was a corpse; “neither eye nor flesh  
 1471 nor any part remains except equivocally”<sup>5</sup> (*QDA*, a. 10 resp.). At the moment  
 1472 that Socrates dies, his body ceases to exist, and a corpse (or, more accurately, a  
 1473 heap of substances) comes into existence. As all his parts go out of existence  
 1474 when his soul ceases to compose his body, no parts of Socrates are found in  
 1475 his corpse. Socrates’s substantial form informs prime matter directly, and the  
 1476 *only* matter that persists over a change of substances is prime matter.

1477 Stump’s characterization of a substance as an “emergent whole” aims to  
 1478 capture this relation by noting that the parts actually composing a substance  
 1479 are not identical with the things that potentially compose it, and similarly  
 1480 for the properties and powers of those parts. A substantial form is precisely  
 1481 that form which accounts for the existence of material substances in general,  
 1482 including those that might be part-less simples, and therefore accounts not  
 1483 only for the composition of some parts into a whole but for the *matter* of the  
 1484 whole. Aquinas draws a distinction between two senses a material composite  
 1485 can have matter. One is the familiar sense in which a material substance has  
 1486 its *integral* parts, such as my hands or fingers or toes, which are the material  
 1487 parts composing me. These are the “proximate matter” of a material substance.  
 1488 Another sense is much less familiar. This is the way of considering matter in  
 1489 a general sense as a part of any material composite substance. And Aquinas,  
 1490 indeed, is known for characterizing this sense of matter as a *potentiality*. This  
 1491 potentiality is what Aquinas calls “prime matter” (*materia prima*):

1492 Prime matter is [...] matter without any form at all, “materiality”  
 1493 (as it were) apart from configuration. When it is a component in  
 1494 a matter-form composite, prime matter is the component of the  
 1495 configured composite which makes it the case that the configured  
 1496 thing can be extended in three dimensions and can occupy a par-  
 1497 ticular place at a particular time. But by itself, apart from form,  
 1498 prime matter exists just potentially; it exists in actuality only as  
 1499 an ingredient in something configured. So we can remove form  
 1500 from prime matter only in thought; everything which exists in  
 1501 reality is configured in some way. For this reason, Aquinas some-  
 1502 times says that form is the actuality of anything. Configuration

---

5 “recedente anima, neque oculus neque caro neque aliqua pars remanet nisi aequivoce.”

1503 or organization is necessary for the existence of anything at all;  
 1504 without form, nothing is actual. (Stump 2003, 37)

1505 Aquinas's prime matter is not one obscure material stuff that is a part of  
 1506 every object, an ultimate material substrate of which everything else is a  
 1507 modification (*pace* Jaworski 2016, 332). Aquinas is strongly against such a  
 1508 theory, in fact, as he argues that prime matter must not be a special kind of  
 1509 matter, in the sense in which my flesh or calcium are kinds of matter, but  
 1510 needs to be devoid of all form. If all material objects had one substratum with  
 1511 its own form, and this substratum was part of every material object, he argues  
 1512 that the substratum would be the only true substance and every other object  
 1513 would be a modification of it (see Thomas Aquinas, *De Substantiis Separatis*,  
 1514 c. 6). It would be the inverse of an atomistic universe, with all everything  
 1515 forming one "bobject."

1516 Prime matter is thus not an integral part but the potential to be a material  
 1517 object, considered apart from any particular actual way something could be  
 1518 a material object by being a member of a determinate kind of thing. The  
 1519 characterization of prime matter as the potentiality for a thing to have a  
 1520 location in space-time and extension in three dimensions follows from the  
 1521 fact that Aquinas holds these features as proper to all material objects in  
 1522 general, of any kind. These features of matter in general are not merely a  
 1523 relation or a feature of our concepts or definitions of matter but are essential  
 1524 to matter in general, i.e., "the potentiality of matter is nothing other than  
 1525 its essence" (*ST*<sup>6</sup> I, q. 77, a. 1, ad. 2; see also Wippel 2000, 319). Aquinas's  
 1526 claims about prime matter are therefore claims that what is essential in being  
 1527 a "material object" is that something has dimensions and spatiotemporal  
 1528 configuration, but not that one has those features in any determinate way.  
 1529 For something to be "material" is only to have indefinite dimensions and  
 1530 space-time location. Prime matter plays a role aside from being the principle  
 1531 in virtue of which things have dimensions because, as we will see, prime  
 1532 matter plays a theoretical role in how we should understand certain kinds of  
 1533 material changes. Prime matter is the matter *from* which some set of integral  
 1534 material parts are constituted (*pace* Dumsday 2021, 2721; see further *ST* I, q.  
 1535 66, resp.).

1536 Inasmuch as prime matter is only the potential to be a material substance,  
 1537 Aquinas therefore holds that prime matter could not exist without having a  
 1538 substantial form to give it particular properties and to determine what *actual*

---

6 Thomas Aquinas, *Summa Theologiae*.

1539 dimensions or location would have. Moreover, Aquinas argues that prime  
 1540 matter alone *necessarily* cannot constitute any particular kind of object, as  
 1541 it is strictly *contradictory* to claim that prime matter exists without being  
 1542 informed by any substantial form. For prime matter to exist by itself, without  
 1543 the actuality of any form, would be to say that something purely potential  
 1544 could be actually existent yet without being actual in any way. That would  
 1545 be nonsense (*ST* I, q. 66, a. 1; cf. *Quodlibet* III, q. 1, a. 1, resp.). All this is to  
 1546 say, in short, that prime matter is a *metaphysical part* of material composites,  
 1547 having a similar kind of relation to a substance as that which a substantial  
 1548 form does: as Aquinas puts it, prime matter is “*incomplete* being without the  
 1549 substantial form” (Marmodoro and Page 2016, 12).

1550 Aquinas’s idea can then be put much more simply: substances are the  
 1551 place where the buck stops for existing, being actual, composing, or having  
 1552 properties or powers. Substances are what really exist, whereas parts only  
 1553 have existence insofar as they are parts of substances: “Only the composite  
 1554 whole [viz., a substance] has existence (*esse*), properly speaking. All the other  
 1555 parts of the substance[, ...] though things of a certain kind, nevertheless exist  
 1556 only in an improper sense, in virtue of the whole’s existence” (Pasnau 2011,  
 1557 624). Aquinas therefore distinguishes two ways of attributing existence to  
 1558 things either as substances or as modifications “in” another:

1559 Existence (*esse*) is attributed to something in two ways. In one  
 1560 way, as to that which properly and truly has existence or exists,  
 1561 and in this way it is attributed only to a substance that subsists *per*  
 1562 *se*. Thus *Physics* I [186b4–8] says that a substance is what truly is.  
 1563 All those things, on the other hand, that do not subsist *per se*, but  
 1564 are in another and with another—whether they are accidents or  
 1565 substantial forms or any sort of parts—do not have existence in  
 1566 such a way that they truly exist, but existence is attributed to them  
 1567 in another way—that is, as that by which something is—just as  
 1568 whiteness is said to be not because it subsists in itself, but because  
 1569 by it something has existence-as-white (*esse album*).<sup>7</sup> (Thomas

---

7 “Uno modo ut sicut ei quod proprie et vere habet esse vel est. Et sic attribuitur soli substantiae per se subsistenti: unde quod vere est, dicitur substantia in I Physic. Omnia vero quae non per se subsistunt, sed in alio et cum alio, sive sint accidentia sive formae substantiales aut quaelibet partes, non habent esse ita ut ipsa vere sint, sed attribuitur eis esse alio modo, idest ut quo aliquid est; sicut albedo dicitur esse, non quia ipsa in se subsistat, sed quia ea aliquid habet esse album.”

1570 Aquinas, Quodlibet IX, q. 2, a. 2, resp.; translated by Pasnau 2011,  
1571 624)

1572 Aquinas, therefore, defines forms in terms of their role, forms being that in  
1573 virtue of which something has actual existence in some way: “all that from  
1574 which something has existence [*esse*], whether that existence is substantial or  
1575 accidental, is able to be called a ‘form’ [...] and because form makes [some-  
1576 thing] to be in actuality, therefore form is said to be an actuality [*actus*]”  
1577 (Thomas Aquinas, *De Principiis Naturae*, caput 1, 5). As prime matter lacks  
1578 any actuality and is essentially a potentiality for being configured by a form,  
1579 the substantial form of a given substance accounts for everything in terms of  
1580 which that material thing is a determinate member of its kind, e.g., having  
1581 essential properties or powers (Stump 2003, 38). According to Aquinas’s way  
1582 of thinking, “there is no such thing as existence beyond the specific ways of  
1583 functioning manifested by specific kinds of things” (Pasnau 2012, 492). This  
1584 is why Aquinas claims that a “substantial form gives being [*esse*] to matter  
1585 *simpliciter*” (SLM,<sup>8</sup> 775).

1586 Further, corresponding to the way in which there are two senses as to what  
1587 the “matter” is in a material substance, Aquinas distinguishes two senses of  
1588 what potentiality in a substance that the substantial form actualizes. Even  
1589 though the matter *from which* every material composite is constituted is prime  
1590 matter, and “no other substantial form intervenes between [a substantial form  
1591 such as] the soul and prime matter” (QDA, a. 9, resp.), no object is merely ac-  
1592 tualized prime matter. Instead, the matter *of which* some material substance is  
1593 composed is its proximate matter, e.g., its integral parts. Hence, the immediate  
1594 potentiality that the human soul makes actual is a living human body (and  
1595 all its parts); “the human body is the matter proportionate to the human soul,  
1596 and it is related to the soul as potency to actuality” (QDA, a. 1, ad 5). Aquinas  
1597 considers such proximate matter brought into existence by the substantial  
1598 form as a *particular way* in which the potentiality of prime matter is actual-  
1599 ized. Since no one substantial form actualizes all of the potential of prime  
1600 matter, this is only some specific potentiality of matter that corresponds to the  
1601 specific actuality that a substantial form brings about: whatever matter, under  
1602 whatever determinate conditions, that is essential to the kind of substance  
1603 the form constitutes.

1604 Consequently, integral parts are *actually* what they are only in virtue of  
1605 composing their substance. The actuality of the parts “is in some sense derived

---

8 Thomas Aquinas, *Sententia libri Metaphysicae*.

1606 from the actuality of the whole, inasmuch as the whole substance, including  
 1607 all of its parts, shares in just a single existence” (QDA, a. 1, ad 5). Conversely,  
 1608 every part of a substance, merely by being a part, is something “in potentiality”  
 1609 (*in potentia*) to the substantial form of that substance. And Aquinas draws  
 1610 this conclusion quite clearly:

1611 That parts [of a substance] are in potentiality alone is apparent  
 1612 because none among them is separate, inasmuch as, given that  
 1613 all the parts, insofar as they are parts, are united in a whole. For  
 1614 everything which exists actually ought to be distinct from other  
 1615 things, because one thing is distinguished from another by its  
 1616 own actuality and form [...]. But those things, which are taken to  
 1617 be parts, are separated from each other when the whole dissolves,  
 1618 then indeed they are beings in actuality, surely not as parts, but  
 1619 as matter existing in privation from the form of the whole. Just as,  
 1620 clearly, in the case of earth, fire, and air, which, when those are  
 1621 parts of a mixed body, are not actually in existence, but only poten-  
 1622 tially [existing] in a mixture. When they become truly separated,  
 1623 then they are in actual existence and are not parts. For, none of  
 1624 the elements, before they are arranged (that is, before they are  
 1625 altered in the mixture and become one mixed thing [composed]  
 1626 from those elements), is one [element] with another, except in  
 1627 the sense that a heap of stones is one thing *secundum quid* [i.e.,  
 1628 in some qualified sense] and not simply.<sup>9</sup> (SLM, 1632–1633)

1629 Obviously, on this way of understanding forms as being that in virtue of which  
 1630 not only some parts are configured into a composite material substance, but  
 1631 that all of that substance’s matter *exists*, it is not easy to see what pluralism  
 1632 about substantial forms could mean. If a substance had two substantial forms,  
 1633 this would be for one and the same substance to exist “twice over,” and that

---

9 “Et quod partes sint in potentia tantum, patet, quia nihil de numero earum est separatum; immo omnes partes in quantum sunt partes, sunt unitae in toto. Omne enim quod est in actu, oportet esse ab aliis distinctum, quia res una dividitur ab alia per suum actum et per formam, sicut supra dictum est. Quando autem ea, quae ponuntur partes, fuerint separata ab invicem dissoluto toto, tunc quidem sunt entia in actu, non quidem ut partes, sed ut materia existens sub privatione formae totius. Sicut patet de terra et igne et aere, quae quando sunt partes corporis mixti, non sunt actu existentia, sed potentia in mixto; cum vero separantur, tunc sunt in actu existentia, et non partes. Nullum enim elementorum antequam digeratur, id est antequam per alterationem debitam veniat ad mixtionem, et fiat unum mixtum ex eis, est unum cum alio, nisi sicut cumulus lapidum est unum secundum quid, et non simpliciter.”



1634 just seems nonsense. Aquinas therefore treats pluralism about substantial  
1635 forms as a conceptual confusion: “Since every form gives a certain *esse*, and  
1636 it is impossible for one thing to have two substantial existences (*esse*), it is  
1637 necessary that if the first substantial form coming to matter gives substantial  
1638 *esse* to it, a second superadded form must give an accidental existence (*esse*)”  
1639 (*In II Sent.*, dist. 18, q. 1, art. 2, corp.; translated by Wippel 2011, 122). Forms  
1640 either make a substance exist, simply speaking, or they otherwise configure  
1641 that substance to exist in some way (e.g., as having a property). And if the  
1642 substance already exists, any further forms in the substance can only bring  
1643 about modifications within that already-existing substance; i.e., further forms  
1644 would not be *substantial* forms.

1645 It is easy to misunderstand Aquinas’s claim about parts as in potential  
1646 to the substantial forms of the substances they compose. Pasnau sees in  
1647 Aquinas’s claim that parts of substances only exist potentially as an attempt  
1648 “to distinguish a thing from its existence, as if it is one kind of question to  
1649 ask whether a thing is real, and another kind of question to ask whether it  
1650 exists” (Pasnau 2011, 627). Pasnau is assuming that Aquinas’s view is that one  
1651 and the same thing becomes a potential thing when it is a part and then is an  
1652 (actually) existent thing when it becomes a substance, where both potential  
1653 and actual things are real.<sup>10</sup> But Aquinas is more radical: he is not going to  
1654 countenance material parts or wholes “surviving” substantial changes of these  
1655 sorts. As Koslicki notes of Aquinas’s views, “no object that is *not* already part  
1656 of a whole that is unified under a single form can survive *becoming* part of  
1657 such a whole; and no object that *is* already part of such a whole can survive  
1658 *ceasing* to be part of it” (Koslicki 2008, 147). So it would be strictly false, on  
1659 Aquinas’s view, that one and the same thing could be characterized at one  
1660 time as a part and at another as a substance. Potential parts of substances are  
1661 not the same things that are the actual substances they can become.

1662 Aquinas, therefore, also treats the view that a substance can have other  
1663 substances as parts as a conceptual confusion. Having other substances as  
1664 proper parts is just what it is to be an aggregate and not a substance. Substantial  
1665 forms, on Aquinas’s view, account for the existence of a substance precisely  
1666 because they account for the existence of every part of that substance:

1667 The soul as the form of the body [...] is united directly to the whole  
1668 body, because it is the form of the body as a whole and of each

---

10 Pasnau thinks Aquinas’s response requires an appeal to the doctrine of a “real distinction” between “essence and existence” (2011, 626–627).

1669 of its parts. And this must be maintained, for, since the body of a  
 1670 man or that of any other animal is a certain natural whole, it will  
 1671 be said to be one because it has one form whereby it is perfected,  
 1672 and not simply because it is an aggregate or a composition, as  
 1673 occurs in the case of a house and other things of this kind. Hence  
 1674 each part of a man and that of an animal must receive its act of  
 1675 existing and species from the soul as its proper form. (*QDA*, a. 10,  
 1676 resp.; translated by John Patrick Rowan in 1949)

1677 Notice, however, that this way of thinking entails that substantial form is  
 1678 *intrinsic* to the substance and all of its parts. An arrangement, for example,  
 1679 is not something intrinsic to the things arranged, and this is what makes an  
 1680 arrangement an accidental rather than a substantial form—it is only a relation  
 1681 among substances. Aquinas uses the illustration that a mass of bronze coming  
 1682 to be a statue only involves an accidental change or alteration because “the  
 1683 bronze, before the advent of the form or figure, has actual existence, and its  
 1684 existence does not depend on that figure” (*De Principiis Naturae*, caput 1, 8).  
 1685 If the statute’s shape were a substantial form, that shape would not only result  
 1686 in the existence of bronze shaped-as-a-statue but the existence of its matter as  
 1687 well. “A form must be something *of* that to which it gives existence, for form  
 1688 and matter are intrinsic principles constituting the essence of a [corporeal]  
 1689 thing” (*De Principiis Naturae*, caput 1, 8).

1690 Similarly, a substantial form is not like a causal *agent* internal to some parts,  
 1691 e.g., gathering them together or pushing them through space. To say then that  
 1692 a substantial form is that in virtue of which a substance exists or is actual is not  
 1693 to say that the substantial form creates or generates its own material parts. A  
 1694 chemist who makes a new chemical compound by combining the constituents  
 1695 in the right way is bringing into existence that compound, certainly, but in  
 1696 a different sense. Aquinas thinks of a causal agent as making some matter  
 1697 *to have a form*: “corporeal forms are caused [...] by matter being brought  
 1698 from potentiality into act by some composite agent” (*ST I*, q. 65, a. 4, resp.;  
 1699 translated by the English Dominican Fathers in 1912–1936). This account of  
 1700 causal agency, even though utilizing an act-potency distinction, presumes  
 1701 that forms play a distinct role. The forms are that *in* what is actual, whereas  
 1702 the agent remains outside of what she actualizes. Whereas the chemist does  
 1703 not become the chemical compound she mixes up, the substantial form and  
 1704 the matter actualized by it “must have one and the same act of existing (*esse*),

1705 something which is not true of an efficient cause and an effect to which it  
 1706 gives esse” (Wippel 2011, 124).

1707 Aquinas’s controversial claim can then be stated more fully as follows: the  
 1708 substantial form not only accounts for the existence of the substance and the  
 1709 composition of other material parts in a whole substance but for *everything that*  
 1710 *is essential to the parts*, whether existence or actuality or powers or properties.  
 1711 As we saw, this claim entails that a substance ceases to exist when it begins to  
 1712 compose a part of something else. Then, given that the substance no longer  
 1713 exists when it becomes a part, all of its properties or powers also cease to  
 1714 exist. Similarly, composing a whole with certain properties, like being human,  
 1715 entails that the parts also have certain properties in virtue of being parts. Thus,  
 1716 my hand is a human hand merely in virtue of composing me, but it ceases to  
 1717 be a hand when it ceases to compose me.

## 1718 2 The Puzzle of Parts

1719 Now, a puzzle looms. Aquinas’s claims about Socrates’s hand ceasing to be a  
 1720 hand when he dies or his body ceasing to be a body both seem empirically  
 1721 false. Consider a case presented by William Jaworski as a counter-example to  
 1722 Aquinas’s theory of composition:

1723 OXYGEN. In the process of respiration, oxygen atoms, as molecular  
 1724 oxygen (O<sub>2</sub>), enter a human bloodstream. Those atoms oxidize red  
 1725 blood cells, becoming parts of those cells and, by extension, a human  
 1726 being. After circulating in human blood, those same oxygen atoms  
 1727 are eventually expelled, albeit in a different molecular configuration  
 1728 (CO<sub>2</sub>).<sup>11</sup>

1729 If we assume that oxygen atoms, molecules, and human beings are all sub-  
 1730 stances, Aquinas is apparently committed to saying that these oxygen atoms  
 1731 were not the same atoms at every point in this process. As Jaworski puts it,  
 1732 “That atom does not survive being incorporated into me. It is instead replaced  
 1733 by something else—something that perhaps has many of the same charac-  
 1734 teristics as the original atom, but that is nevertheless numerically different  
 1735 from it” (2016, 118). Aquinas appears, therefore, to claim that when those  
 1736 atoms begin to compose a human being, those atoms ipso facto cease to exist.

---

11 I have made the case more specific. See below.

1737 In fact, Aquinas would be committed to the stronger claim that those atoms  
 1738 *never* existed because they always composed some other substance at every  
 1739 step of the case.

1740 Yet we see no such replacements happening when substances come to  
 1741 acquire new parts; the oxygen atom does not appear to be replaced by a  
 1742 “token” oxygen look-alike. We could have used radioactive isotopes to “tag”  
 1743 the atoms within the molecular oxygen and then identify the same two atoms  
 1744 at every point in the process. If these atoms ceased to exist or never existed,  
 1745 how could we track each particular atom, their properties, and their causal  
 1746 powers? Oxygen atoms do not disappear when they compose other molecules,  
 1747 nor do their properties or powers just cease to exist. Cases like OXYGEN are  
 1748 not exceptional or infrequent. When we break up a composite substance, the  
 1749 ingredient substances can come back into full existence with entirely the same  
 1750 properties they had before they composed anything. Oxygen atoms do not just  
 1751 “pop” into existence when we break up, e.g., H<sub>2</sub>O molecules with hydrolysis;  
 1752 the atoms were parts of the molecular structure itself! The Thomist view thus  
 1753 appears straightforwardly empirically false.

1754 This puzzle should not be as puzzling as it might seem. Aquinas holds that  
 1755 substances cannot be parts of other substances, and if an atom becomes a part  
 1756 of a molecule, ipso facto, that atom ceases to be a substance. Nevertheless,  
 1757 Aquinas does not hold merely that the atom no longer exists. Rather, his claim  
 1758 is simply that the thing that was an atom substance *became* an atomic part of  
 1759 a molecule. For Aquinas, if the oxygen atom was a substance and remained a  
 1760 substance over the event described in OXYGEN, the atom would not compose  
 1761 that molecule but only become, at best, spatially co-located with the molecule.  
 1762 More accurately, as molecules are not separable things from the atoms that  
 1763 compose them, it would be that molecules are nothing more than spatial  
 1764 arrangements of atoms; i.e., molecules are not genuine material objects but  
 1765 pseudo-objects.

1766 Further, every substance exists in a determinate way, that is, as a member  
 1767 of a kind. If *oxygen atom* and *hemoglobin protein* are each a distinct kind, it  
 1768 is easy to see that the implicit assumption is likely empirically false. When  
 1769 incorporated into an oxyhemoglobin molecule, an oxygen atom is config-  
 1770 ured differently from when it is not incorporated into a protein. When, in  
 1771 respiration, oxygen atoms are incorporated into the hemoglobin in red blood  
 1772 cells, those atoms bond with the hemoglobin, and their structure changes. A  
 1773 free-floating oxygen atom undergoes a series of changes when it bonds with  
 1774 hemoglobin, such that it comes to have different properties and structural

1775 relations to other things (e.g., the hemoglobin—see, for example, Jenkins and  
1776 Van Kessel 2002, 122). In becoming a part of a protein, then, the Thomistic  
1777 account holds that the oxygen atom ceases to be a substance when it becomes  
1778 a structural component of the substance that is the protein. It is no part of the  
1779 account, *pace* Jaworski’s objection, that the oxygen atom is replaced with a  
1780 completely identical look-alike when it bonds with hemoglobin; instead, the  
1781 oxygen atom becomes quite different structurally and in its other properties  
1782 at the moment it becomes a structural component of a protein.

1783 Consider for a moment the simpler case of H<sub>2</sub>O molecules. H<sub>2</sub>O and O<sub>2</sub> are  
1784 distinct molecules with distinct properties and powers. These two molecules  
1785 have distinct properties and powers because they are distinct structurally.  
1786 Further, their oxygen parts are distinct structurally as well: the oxygen in  
1787 H<sub>2</sub>O has two distinct covalent bonds with hydrogen, and dioxygen’s atom  
1788 parts have a double covalent bond. These kinds of bonds modify the distinct  
1789 properties of the whole in a way that the whole has properties and powers  
1790 distinct from other possible configurations, but it is also true that the oxygen  
1791 atoms being so bonded are distinct in powers and properties from a single  
1792 free-floating oxygen atom. Thus, an individual atom might react under certain  
1793 conditions (hydrogen gas will react with O<sub>2</sub> in combustion), whereas in the  
1794 molecule it does not react so (H<sub>2</sub>O does not combust, even as a gas). If each  
1795 constitutes a distinct kind of substance, it is not clear how molecular oxygen  
1796 coming to compose a thing of a distinct kind has *not* ceased to be an instance  
1797 of “molecular oxygen.”

1798 The only way in which the case of OXYGEN could be a counter-example to  
1799 Aquinas’s theory of composition is if it described a case where something came  
1800 to compose another without any change of essential properties. First, the case  
1801 does not plausibly show this, since it is an empirical matter whether there *is*  
1802 such a natural kind as an “oxygen atom,”<sup>12</sup> and the atom in the case underwent  
1803 a great deal of changes that contrast with the way that we ordinarily take an  
1804 atom to be determinate when it is *not* part of those compounds. Second, it is  
1805 not clear how such a case *could* disprove Aquinas’s views without assuming  
1806 what it intends to refute. Aquinas’s views are that something only counts as  
1807 an instance of material composition when the parts depend on the whole in  
1808 a certain way. If the case was taken to describe merely extrinsic changes of  
1809 spatial location among the atoms, Aquinas would just flatly deny that the

---

12 Kerry McKenzie has written extensively in criticism of the view that particles are fundamental entities; see McKenzie (2014, 2011) and McKenzie and Muller (2017).

1810 atoms composed anything throughout the process. Instead, Aquinas's claim  
1811 only entails that the atom, when it composes a molecular substance, is at  
1812 best a derivative property-bearer in virtue of that substance (if it is the right  
1813 integral part of the molecule to do so) and that any properties it has that result  
1814 from being a part of a molecule would cease when it ceases to compose that  
1815 molecule.

1816 In fact, the puzzle of parts is actually not a puzzle about material composi-  
1817 tion but a puzzle about what persists over substantial change or the change of  
1818 one substance into another substance (there would be *no* puzzle if there were  
1819 no changes of parts). And Aquinas's perspective, bluntly, is that the wrong  
1820 place to look for continuity in change is in what it is to be a substance. The  
1821 right place to account for continuity in substantial change is in the relation  
1822 among the substances that go into or come out of existence while carefully  
1823 distinguishing the material parts involved in the changes. The reason that  
1824 H<sub>2</sub>O can be split into oxygen and hydrogen by electrolysis is not that there  
1825 was both molecular oxygen and a water molecule spatially co-located at the  
1826 beginning of the process. Rather, the reason one can split out these two com-  
1827 ponents is because water molecules are such that they can be decomposed  
1828 into hydrogen and oxygen atoms.

1829 What we should appeal to in order to account for continuity in substantial  
1830 changes is the actuality or potentiality corresponding to the substance (and its  
1831 integral parts') ability to undergo the relevant changes. Hydrogen and oxygen  
1832 atoms are "potential parts" of water molecules because they are the proximate  
1833 matter essential to being a thing of the kind "water molecule." H<sub>2</sub>O does not  
1834 exist without them. But it is not essential to hydrogen or oxygen atoms to  
1835 constitute H<sub>2</sub>O. Further, we can assume H<sub>2</sub>O molecules have distinguishable  
1836 integral parts, such that we can identify the hydrogen and oxygen atomic  
1837 parts; that is, the parts can have their own properties that they bear in virtue  
1838 of being parts of that substance (the parts bear properties derivatively). Then,  
1839 we can say that when the water molecule is decomposed in hydrolysis, there  
1840 are two senses in which *the same* integral parts are what became one hydrogen  
1841 and two oxygen atom substances.

1842 On the one hand, the matter *from which* they were constituted is just the  
1843 same matter that came from the molecule because their coming into exist-  
1844 ence consists in an actualization of that potential—prime matter—that was  
1845 formerly "in" the water molecule. They did not "pop" into existence from  
1846 nowhere. On the other hand, the material integral parts of the water molecule  
1847 were characterized as (derivative) property bearers that had their own inter-

1848 nal structures and properties in virtue of being essential parts of the water  
1849 molecule. When the water molecule decomposes into the atoms, the new sub-  
1850 stances only need to lose those properties that were essential to the whole they  
1851 composed. We can imagine, for example, that the hydrogen atomic part being  
1852 tagged with an isotope is a property that is accidental to both the molecule  
1853 and the hydrogen. If they bore any properties that were unique to themselves  
1854 as the parts, these would be accidental to the water molecule and could come  
1855 to characterize the new substances as well. There is nothing preventing the  
1856 hydrogen atom substance, resulting from the decomposition of the molecule,  
1857 from likewise being characterized by the isotope tagging.

1858 It might be alleged that, in these cases, “scientists do not claim to be tracing  
1859 powers, but things that bear the powers” (Pawl and Spencer 2016, 138). But  
1860 Aquinas’s view does not require thinking that all we track are only powers of  
1861 oxyhemoglobin. Instead, some parts can bear properties insofar as a substance  
1862 can have a property *in*, or in virtue of, one of its parts,<sup>13</sup> where the whole will  
1863 be the subject even of accidental properties of the parts, as when, for example,  
1864 I have the property of “being wounded” in virtue of my foot being wounded.  
1865 And one of these integral parts can be such that it can become a numerically  
1866 distinct substance, bearing that property in its own right, when it ceases to  
1867 compose the whole. Atomic parts are just like this. There does not seem to  
1868 be any empirical reason to think that, in tracking an oxygen atom through  
1869 my body with a radioactive isotope, we need more than Aquinas’s account  
1870 can give: a certain isotope was introduced into my body, and, in virtue of a  
1871 chemical change, became composed as a part of one of my atomic parts; that  
1872 atomic integral part of me was tracked, in virtue of the radioactive properties  
1873 now associated with that atomic part, and then the isotope part or the atomic  
1874 part ceased to compose me, eventually (see further Toner 2008).

1875 Yet, Pasnau alleges there is the *inverse* problem of that posed by Jaworski:  
1876 how to explain the fact that *exactly similar* properties persist over substantial  
1877 changes. His example is that the skin color of Socrates can be identical to  
1878 the skin color of Socrates’s corpse a moment after death, and “it seems noth-  
1879 ing short of miraculous that, without that form, the corpse retains so many  
1880 exactly similar accidents” (Pasnau 2011, 585). Here again, it is important to  
1881 note that while it is true that Socrates’s substantial form is that in virtue of  
1882 which Socrates and his parts are characterized by essentially human proper-

---

13 Aquinas draws this very contrast between properties and parts. Properties, accidents, are not “particular things.” But parts, even though they too are dependent entities like properties, can be considered particular things in ways that properties cannot (QDA, a. 1, ad. 9).

1883 ties, Aquinas’s claim is *not* that no qualitatively similar kind of property or  
 1884 part, even an *exactly* similar property, could ever characterize anything else  
 1885 (Quodlibet I, q. 4, a. 1; translated by Sandra Edwards in 1983). There is no  
 1886 reason that Socrates and Socrates’s corpse could not be qualitatively identical  
 1887 in regard to skin color (and Aquinas does say that they could be; *De Ente et*  
 1888 *Essentia*, c. 5). Aquinas’s view only requires holding that Socrates’s corpse  
 1889 is not the same substance as Socrates’s body and that no essentially human  
 1890 properties “survive” Socrates’s demise.

1891 Some kinds of substances, given their proximate matter, have the potential  
 1892 to become other substances, whereas others have integral parts such that  
 1893 those parts can become substances in various ways. For that reason, organ  
 1894 transplants are not a metaphysical mystery. A heart, when detached and “on  
 1895 ice,” is no longer a part of any particular human, although it is suited to become  
 1896 the heart of another person because of its physical characteristics; the heart  
 1897 can retain those properties, while detached, that did not derive solely from  
 1898 composing a human being. E.g., muscular cells are still capable of moving  
 1899 under electric shocks, and the whole heart is capable, when reattached, of  
 1900 pumping blood (*pace Pawl and Spencer 2016, 144*). Nothing about Aquinas’s  
 1901 position requires that the heart, when it is *in via* during a transplant, will not  
 1902 be a thing “suited to beat and pump blood.”<sup>14</sup> All that is required, on Aquinas’s  
 1903 metaphysics, is that my heart has undergone some intrinsic, essential change  
 1904 when it ceased to compose my body, such that it is a distinct thing when it is a  
 1905 part of me and when it is not. What we want to know is why, if it is a distinct  
 1906 thing, the heart outside of my body has apparently very similar properties. On  
 1907 Aquinas’s view, the answer is that my heart was just the kind of part that could  
 1908 become such a substance—acquire *that* kind of substantial form—when it  
 1909 was detached from my body, given the proximate matter of which human  
 1910 being from which it was taken was composed essentially included a heart.

1911 In **OXYGEN**, the relation among the substance kinds to which the oxygen,  
 1912 hemoglobin, and oxyhemoglobin belong explains the potentiality of the oxy-  
 1913 gen atom to become a part of oxyhemoglobin. Other things could not compose  
 1914 hemoglobin unless they both underwent some suitable external stimulus *and*  
 1915 were suitable to have the potentialities to compose oxyhemoglobin. However,  
 1916 in fact, all of these conditions were met in **OXYGEN**. Thus, when actual-  
 1917 ized by some stimulus conditions, the proximity of the oxygen atom to the

---

14 I am speaking generally because a detached heart is likely not a substance but a collection of individual substances (cells).



1918 hemoglobin initiated a chemical reaction of oxidization of that protein, and  
1919 oxyhemoglobin was composed from those other substances. Similarly, after it  
1920 comes to composing oxyhemoglobin, that oxygen atom substance becomes  
1921 an atomic part typical of oxyhemoglobin and has a certain set of chemical  
1922 bonds with the protein. Not every oxygen part of any molecule is of such a  
1923 type as *this* oxygen part of oxyhemoglobin: not every oxygen part bonds with  
1924 the particular geometry involved in an oxygen part's bond to the rest of the  
1925 oxyhemoglobin molecule (i.e., an "end-on bent" configuration in bonding  
1926 with the Fe<sub>2</sub> parts of that molecule).<sup>15</sup> Yet some other kinds of molecules  
1927 could have oxygen in the same type of configuration as the oxygen parts of  
1928 oxyhemoglobin.

1929 An objector might point out that the homonymy principle entailed, for  
1930 Socrates, that Socrates's eye is no longer an eye after he dies. The objector could  
1931 then argue that we should not think that "atom" is being used homonymously  
1932 of the atom substance and the atomic part of oxyhemoglobin: "These are  
1933 'atoms' in just the same sense, whether or not they compose anything! Whereas  
1934 it might be plausible that 'eye' is a functional term for a certain kind of part,  
1935 and we can imagine it ceases to apply to an eye when it is separated from its  
1936 functional system, surely atoms are not a functional part of that sort."

1937 In response, first, it seems likely to me that Aquinas and Aristotle treat the  
1938 aforementioned "transplant" cases as the organ ceasing to have any biological  
1939 properties merely because organ transplants were not then medically possible,  
1940 and they did not know that an organ's cells do not immediately cease to  
1941 be alive on detachment. Yet, even if we *were* committed to the homonymy  
1942 principle for all parts, this can be plausible when we specify the nature of  
1943 the kinds in question. If we assume that kinds are kinds of *substances*, and  
1944 substances are those objects composing no other, then a kind such as "oxygen  
1945 atom" cannot apply to the oxygen in oxyhemoglobin. An oxygen atom as a  
1946 substance is, by stipulation, something that does not compose anything else,  
1947 and the oxygen in hemoglobin clearly composes it. As an integral part that  
1948 essentially characterizes molecules of the kind, the atomic part now belongs  
1949 to the kind "oxyhemoglobin" in virtue of composing the whole.

1950 The only thing further the objector might be looking for, as we saw with  
1951 Jaworski, is the *numerical identity* of the thing having the property at every  
1952 time it has the property (whether as a part or a whole). But numerical identity  
1953 strikes me as something we cannot just *see*, because we characterize the

---

15 See the case study of carbon monoxide poisoning in Gaffney and Marley (2018, 233–234).

1954 substances (even for Jaworski) in terms of what is essential to them. For  
 1955 something to be numerically identical is to say that it underwent no change in  
 1956 what it is essentially. To say that the atom is “numerically identical” whether  
 1957 it composes the molecule or not entails both that the oxygen differs in *no* way  
 1958 when it composes oxyhemoglobin and that oxyhemoglobin is not a kind of  
 1959 molecule. Both seem empirically false.

1960 The atoms in oxyhemoglobin have distinct shapes, properties, and powers  
 1961 from the oxygen atoms composing O<sub>2</sub>. In addition, if oxyhemoglobin were not  
 1962 a kind of molecule, then the example could not undermine Aquinas’s overall  
 1963 thesis, as the atom only gets spatially located very close to the other. Similarly,  
 1964 if the objector were to insist that numerically identical *properties* characterize  
 1965 the oxyhemoglobin and the oxygen atom that results from its decomposition,  
 1966 even though it is conceded that they are exactly similar, such a response would  
 1967 appear to beg the question against Aquinas that only substances (e.g., atoms  
 1968 and molecules) bear properties. The properties, to be “numerically identical,”  
 1969 would have to be substances, in Aquinas’s sense. If they were substances, how-  
 1970 ever, they could not compose a material object without, necessarily, ceasing  
 1971 to be substances. Thus, it is not clear how to make sense of either attempt to  
 1972 cash out “numerical identity” in a way that does not beg the question against  
 1973 Aquinas’s position.

### 1974 **3 The Plausibility of the Solution**

1975 The ontologically relevant payoff of distinguishing between prime matter and  
 1976 the proximate matter of material parts is that it allows Aquinas to draw a  
 1977 distinction such that he can affirm both that, even though these parts potentially  
 1978 could constitute something else, these parts actually compose a substance’s  
 1979 essential parts. On the one hand, distinguishing a substantial form as a partic-  
 1980 ular, a metaphysical part of a composite, is to say that the substantial form  
 1981 is not identical with those material parts or the whole they compose (*pace*  
 1982 [Scaltsas 1994](#); cf. *SLM*, 1674, translated by John P. Rowan in [1961b](#); compare:  
 1983 [Keinänen and Hakkarainen 2017](#); [Keinänen 2018](#)). Substantial forms are *par-*  
 1984 *ticulars* that, in virtue of characterizing some set of material parts, account  
 1985 for why those parts constitute a whole of some kind. It is not a feature of  
 1986 our counting or conceptual schemes that the material *x*s are such that they  
 1987 compose *y*, but an extra-mental fact about the *x*s that they compose *y* since  
 1988 they do so by reason of the substantial form that is intrinsic to that material

1989 object and all its parts.<sup>16</sup> This is what it is to say that the substantial form is  
 1990 the actuality of that substance *y* and its parts, the *x*s, is that each of the *x*s  
 1991 (and the substantial form) are such that they *actually* compose *y* in virtue of  
 1992 something that is essential to them (the substantial form<sup>17</sup>), even though the  
 1993 substantial form is not identical with the *x*s or the *y*.

1994 On the other hand, there remains a sense in which *those* material things,  
 1995 the *x*s, could have composed *z* instead of *y*. Insofar as the *x*s are adequately  
 1996 characterized by being “material” (i.e., composed of prime matter as a meta-  
 1997 physical part), Aquinas holds that material objects, in general, are essentially  
 1998 such that they can undergo a change of kind—one material object can serve  
 1999 as matter *from which* we generate the matter of another material object of a  
 2000 different kind, because prime matter is just the potentiality of any material  
 2001 object to come to constitute a distinct kind of material object under the relevant  
 2002 conditions (Brown 2005, 79–83). It is not essential to their matter that  
 2003 the *x*s compose *y*, and thus *z*’s parts could be composed *from* the *x*s.

2004 Aquinas does not hold that matter is fundamentally or essentially particulate,  
 2005 and it is apparent now why he cannot think it is. If matter were  
 2006 fundamentally particulate, it would be essential to those particles that they  
 2007 are mereologically simple, and it would consequently be false that they could  
 2008 ever compose a whole object. Aquinas’s claim that matter is not essentially  
 2009 particulate does not merely constitute medieval empirical speculation lacking  
 2010 knowledge of the existence of fundamental physical particles but follows  
 2011 directly from the assumption that there are material composites (and that  
 2012 any arbitrary two or more material things do not compose an object)—that is,  
 2013 that material composition occurs only under some restricted circumstances.

2014 An account of material composition would be involved in an infinite regress  
 2015 if it only specifies the conditions under which some things come to compose a  
 2016 whole without explaining what the things are that *get* composed. This would  
 2017 be akin to explaining what it is to be a bearer of properties, a thing defined in  
 2018 terms of being what bears properties, by appeal to a distinct property of that  
 2019 thing (see the discussion in Loux 2006, 84–120). Aquinas’s earlier objection  
 2020 to a plurality of substantial forms in one substance is that it involves one in  
 2021 an infinite regress of exactly the same sort; it infinitely defers the question  
 2022 as to what has the potentiality to be composed, or it assumes an ultimately

16 “An act is in that of which it is the actuality: wherefore the soul must exist in the whole body and in each of its parts” (*ST* I, q. 76, a. 8, resp.).

17 By reflexivity, every part is a part of itself. The form is thus a metaphysical part of itself, itself that in virtue of which it composes the whole as a metaphysical part.

2023 atomistic account of reality as a necessary truth (viz., Peter van Inwagen’s  
2024 account of composition—compare: [Renz 2018](#)). If matter were essentially  
2025 particulate, then this would seem to mean that *to be a material object* is just  
2026 to be one of the particles; that is, it is *necessary* for the material objects that  
2027 they be mereologically simple. But this position would involve a confusion  
2028 between two different senses of what it is to be a material object: being the  
2029 sort of thing that essentially has integral parts and being the sort of thing that  
2030 essentially has dimensions, occupying a spatiotemporal location. But it is false  
2031 that having the sort of thing with dimensions and a spatiotemporal location  
2032 necessarily entails that all the material objects are essentially mereologically  
2033 simple.

2034 To put it another way, even though prime matter characterizes every mater-  
2035 ial object, it is only a way to describe the potentiality to be a material object  
2036 and has *no* essential characteristics at all. Prime matter is thus an explanatory  
2037 principle in virtue of which it is a contingent matter whether *any* material  
2038 substance exists—i.e., it is not essential to any material thing, merely in virtue  
2039 of being material, that it be actual. Consequently, whatever constitutes “mater-  
2040 iality” cannot be something that has any essential properties (neither a  
2041 property nor a property-bearer), but rather that special sort of potentiality  
2042 that corresponds to the potential to be a material substance: namely, that no  
2043 material object exists necessarily, but only contingently. Prime matter has to  
2044 be “pure potentiality” in this way in order to thread the needle between the  
2045 views that composition among two or more material things occurs of necessity  
2046 (universalism) and that material objects are essentially mereologically simple  
2047 (nihilism).

2048 In fact, Aquinas argues that if there are things having dimensions and spa-  
2049 tiotemporal location, they are *by that very fact* composite objects—composites  
2050 precisely inasmuch as they possess spatial parts: “From the fact that matter  
2051 has corporeal existence through forms, it immediately follows that there are  
2052 dimensions in matter whereby it is understood to be divisible into different  
2053 parts, so that it can receive different forms corresponding to its different parts”  
2054 (*QDA*, a. 9, resp.; translated by John Patrick Rowan in 1949). Inasmuch as  
2055 material objects have parts that are spatially distinguished, these are *integral*  
2056 *parts*, and we find that distinct integral parts can bear distinct properties or  
2057 evince different structures (that is, spatial parts of one object can bear distinct  
2058 accidental forms). Aquinas’s point is that spatial parts of a substance are the  
2059 right kind of thing themselves to have properties in various ways. We can, for

2060 example, characterize one spatial region of the same substance as hot and  
2061 another as cold because spatial parts can have distinct properties.

2062 John Heil denies this and argues, to the contrary, that if substances can  
2063 only have spatial or temporal parts, then that is enough to claim that they are  
2064 “mereologically simple” (2012, 18–19). The only realistic candidate substances  
2065 would be particles or fields, whereas macroscopic entities like humans are not  
2066 (Heil 2012, 19–22; see also 2003, 177–192). This is because fields or particles  
2067 would not have parts that bear distinct properties from the whole. The whole  
2068 field has one set of properties borne directly by the field, and all of its spatial  
2069 parts (considering the field’s extension in space as a way to divide it into spatial  
2070 parts) have the same properties. Whereas a substance like a field or a particle  
2071 can have many properties, properties are not parts of those things. Properties  
2072 are not parts of their substances (Heil 2012, 107). The argument is to the  
2073 effect that if Aquinas admits that a material object has integral parts, and  
2074 these parts can bear properties, then those parts must be substances. Aquinas  
2075 would therefore be contradicting the shared assumption that only substances  
2076 are property-bearers.

2077 But Aquinas has not assumed that integral parts of a substance, among  
2078 which are that object’s spatial parts, are bearers of properties *in their own*  
2079 *right*. Aquinas just denies this implicit premise of Heil’s argument. What is  
2080 required is a distinction between the fact that some things essentially bear  
2081 properties and that other things bear properties in a derivative way. That is,  
2082 there is no contradiction if integral parts bear properties only *accidentally*, i.e.,  
2083 only in virtue of composing something that is *essentially* a property-bearer, a  
2084 substance. When my hand is white, then *I* am white with respect to my hand.  
2085 My hand is not a property-bearer in its own right but bears properties only  
2086 in virtue of being a part of me. Yet, understanding integral parthood in this  
2087 way, such that integral parts bear properties in an only accidental way, we not  
2088 only can divide an object according to its spatial dimensions—top half, right  
2089 half, etc.—but also in terms of the way in which each distinct part can bear  
2090 distinct properties, or have a power, or be structured.

#### 2091 **4 Conclusion**

2092 The reason that Heil restricts the parts of objects to merely spatiotemporal  
2093 parts, however, seems independently motivated. At root, the difference lies in  
2094 how to identify or classify which things are genuine material objects and, so,  
2095 which things are genuine parts or properties of them. Stump’s elaboration of

2096 the phenomena associated with “emergence” in contemporary metaphysics  
 2097 aims to make it plausible, with appeal to empirical data, that there exist  
 2098 properties that come to be seated in one whole rather than merely a collection  
 2099 of substances. These properties are not such that they could be merely the  
 2100 properties of a complex of substances but must be of an emergent whole. By  
 2101 contrast, Heil is inclined to hold that the scientific data shows that the world  
 2102 is perhaps fundamentally composed only of fields or that they are the only  
 2103 things that qualify, empirically, as having properties.

2104 Further, Aquinas’s substances are emergent wholes in the sense that the  
 2105 whole is not identical with the ingredients that lead to their emergence because  
 2106 he believes substantial change of one substance into another is possible.<sup>18</sup>  
 2107 Heil would likely disagree with the terminology of “emergence” because, e.g.,  
 2108 it is not clear under what circumstances two fields could come to compose a  
 2109 distinct object. One field does not appear to be the thing that could become  
 2110 another field. Heil’s vision then seems to entail that there is only accidental  
 2111 change among the fields that exist.

2112 Yet Aquinas is not claiming that it is necessarily the case that there are  
 2113 material objects, but only that there would be no problem of material com-  
 2114 position if there were no material objects. From Aquinas’s perspective, if the  
 2115 scientific data shows us that the world is built entirely from fields, and fields  
 2116 are things that entirely lack spatial parts and spatiotemporal location, then  
 2117 it would merely be the case that the existence of spatiotemporally-located  
 2118 objects is merely *apparent*. (Of course, one would need to explain how the

---

18 Even if Aquinas were wrong about the substantial change of macroscopic entities, the claim becomes far more plausible in the subatomic world. The identities of some physical particles appear inseparable from the physical systems they form—for example, when electrons become “entangled” in a quantum state. Aquinas could hold that these states are, in fact, hylomorphic composite substances without any parts other than spatial and temporal parts. Electrons appear to cease to exist (except “virtually”) in these states they compose. As with the way Aquinas elsewhere treats an Aristotelian homogenous “mixture”—cf. *De Mixtione Elementorum*—electrons still relate to the subsequent entangled state in virtue of a mathematical “structural” correspondence between the individual quantum state (and powers/properties) of the electrons before their being entangled and that state after they are entangled. This way of considering particles allows us to hold that particles are substances in certain circumstances, even if these particles can come to compose entangled quantum states. This would entail realism about the quantum wave function. The relation of particles as substances to the quantum states would be composition, and of the quantum state to the distinct particles that result from a “collapse” of the wave function would be decomposition. The issues here are obviously highly simplified. See further Ney and Albert (2013). And it seems to me that there are philosophical accounts of the metaphysics of the wave function already being proposed that are compatible with my loose characterization: cf. Gao (2017). Robert Koons (2018) appeals to Aquinas in this way.

2119 fields exist and how these fields generate the apparent spatiotemporal world.)  
2120 The point is that once we admit that there are objects that have spatial dimen-  
2121 sions, we are committed to the fact that there exist material composites, and,  
2122 subsequently, we can pose questions as to the way in which their composite  
2123 parts compose one whole. Consequently, if there are material objects with  
2124 parts, they at least bear the property of being a *y* such that the *x*s compose it,  
2125 and so these are all ipso facto going to be the kind of thing that can bear (one  
2126 or more) properties by reason of the kind of thing they are.

2127 But this points to an advantage of Aquinas's account of substances. What  
2128 fills out the account, beyond the claims he makes about matter and form  
2129 composition, is the way in which Aquinas thinks we discover and identify  
2130 the natural kinds. Picking out the things with substantial forms thus requires  
2131 us to identify *the causal process* by which some things are modified in order  
2132 to become a whole. There are real differences among the causal processes  
2133 that might produce distinct kinds of substances or parts. The account, in  
2134 fact, rests the nature of what resolves the problem of material composition  
2135 on the extra-mereological considerations that should lead us to hold that  
2136 some causal processes produce (or lead to the destruction of) instances of a  
2137 natural kind. A causal process produces a substance, for instance, when that  
2138 process brings some things together as parts, such that the things that are the  
2139 parts of the resultant substance cannot be described as essentially the same  
2140 as the things that became those parts. Conversely, the process that leads to  
2141 the destruction of a substance is one in which that substance "loses" some  
2142 essential part or property over the course of the change and, ipso facto, ceases  
2143 to exist. Even a "brute" theory of composition on which there is an infinite  
2144 series of conditions under which the *x*s compose *y* can be hylomorphic if  
2145 those conditions are interpreted as corresponding to different natural kinds of  
2146 object; there would be no conflict with hylomorphism's account of the nature  
2147 of material composition if there were infinitely many natural kinds, even if  
2148 there might be some other good reason to think it is impossible for there to  
2149 be truly be an "infinite" number of such kinds (Markosian 2008, 19). The  
2150 account of natural kinds, therefore, tells us which are the things that have  
2151 substantial forms and, so, are the "genuine material objects."

2152 This Thomistic way of approaching material composition makes an empiri-  
2153 cal, scientific method for identifying the substances by appeal to facts about  
2154 causal processes a promising one. What it *means* for a thing to be unified is  
2155 to be a member of a natural kind, to be a material object of a determinate

2156 nature.<sup>19</sup> Whether an atom comes to compose a molecule, for example, is  
 2157 determined by whether that atom's essential properties changed in virtue of  
 2158 becoming an apparent part. These give us those criteria in virtue of which  
 2159 there is unity among the molecule's atomic elements, quarks, and electrons.  
 2160 And those criteria for molecules are of a very different sort from the criteria for  
 2161 other substances. Different natural kinds exemplify different kinds of unified  
 2162 causal powers or activities among the parts; it would be impossible to give a  
 2163 general characterization of what that unity consists in, in general, without  
 2164 a minimal account of the natural kinds. Different natural kinds exemplify  
 2165 different essential activities or properties or powers or structures, and so their  
 2166 parts are unified according to different criteria.

2167 In the Thomistic picture, then, the role that the form has in unifying some  
 2168 parts into a whole depends on *the right kinds of changes in virtue of which some*  
 2169 *object changes kind membership* and thereby becomes a part of another thing  
 2170 of a distinct kind. That is, if substances can never compose other substances  
 2171 as parts, things must cease to be members of one kind and become members  
 2172 of another kind at the time they become parts.<sup>20</sup> For human beings, the  
 2173 unity among our parts is explained by facts about our organic chemistry as  
 2174 animals of a special sort. Captain Hook, as a human, is a living, organic  
 2175 thing of a particular natural kind *human being*, where his parts are unified  
 2176 by biochemical bonds and various bio-organizational interactions. That in  
 2177 virtue of which all of his material parts are of the same kind, *human being*, is  
 2178 what makes those parts belong to Hook, but Hook is neither identical with  
 2179 his kind (the essence of *human being*) nor is Hook merely that which makes  
 2180 him a member of the kind or all his parts human (his soul). Rather, Hook is  
 2181 a substance formed of his material parts, among which is his hand, suitably  
 2182 informed and united by his soul. Captain Hook loses his hand to the crocodile  
 2183 *because* that action causes his hand to cease to have the right biochemical  
 2184 bonds and interactions so that Hook's substantial form ceases to be something  
 2185 Hook shares in common with that piece of matter that formerly was his hand.  
 2186 When Captain Hook acquires a hook made of iron and wood, and puts the  
 2187 hook in the place of his hand, that hook is not a part of Captain Hook because  
 2188 the hook is not changed into the right *kind* of thing that could be a part of his

19 One should note the way this account is similar to the way in which Jansen criticizes and modifies Kit Fine's account of embodiments around the notion of a sortal: see Jansen (2019).

20 Koslicki recognized this implication, where she notes that accepting Aristotle's homonymy principle requires that, if a substance becomes a part, "any such transformation would essentially involve a change in *kind membership*" (Koslicki 2008, 147).



2189 living, organic body. His substantial form cannot overlap the hook because the  
 2190 substantial form just is that in virtue of which Hook is a biological organism  
 2191 of a kind, and the hook is not the right kind of thing to be part of a biological  
 2192 organism.<sup>21</sup>

2193 In sum, on the Thomistic theory, the only way to determine the way in  
 2194 which a substantial form accounts for the unity of the parts of a substance is to  
 2195 determine the nature of the substance in question. This seems the best reason  
 2196 to commend the Thomistic account of substantial forms. Questions about the  
 2197 unity of material objects can be resolved to a certain degree at an abstract,  
 2198 metaphysical level but are fundamentally a matter to be resolved through  
 2199 empirical investigation. Hylomorphism of the Thomistic sort appeals to forms  
 2200 to explain material composition, but what and how a form accounts for the  
 2201 composition of a substance depends on the *kind* of substance it informs. This  
 2202 reliance on natural kinds grounds answers to “what it is in virtue of which  
 2203 the *xs* compose *y*” soundly on empirical concerns; the question of what kinds  
 2204 there are, or their properties, can be given a fully satisfactory answer only in  
 2205 tandem with scientific investigation.<sup>22,23,\*</sup>

2206 Fr. James Dominic Rooney, OP

2207  0000-0003-0087-3218

2208 Hong Kong Baptist University, Hong Kong SAR

2209 jdrooney@hkbu.edu.hk

## 2210 References

2211 BROWN, Christopher M. 2005. *Aquinas and the Ship of Theseus: Solving Puzzles About*  
 2212 *Material Objects*. London: Continuum International Publishing Group.

---

21 For example, a human being *could* digest iron and make it part of its organism, but the hook as-is has not been digested or appropriately modified to form part of Hook’s body. The same would be true of any sophisticated prosthetic; as long as these are such that they are not “biologically continuous” with the human organism, they are not parts of that organism except in perhaps an extended or metaphorical sense. By contrast, if we created a biological replica of Hook’s hand, growing a cloned set of Hook’s cells and structuring them in an appropriate way, then attaching such a biological prosthetic *would* be able to become part of Hook’s body (if it was not rejected by his immune system, etc.).

22 Permission has been given to reprint material adapted from James Rooney (2022).

23 All Latin editions obtained from those hosted at [Corpus Thomisticum](https://www.corpus-thomisticum.com/), courtesy of the University of Navarre.

\* THANKS

- 2213 DUMSDAY, Travis. 2021. "Can a Relational Substance Ontology Be Hylomorphic?"  
 2214 *Synthese* 198(suppl. 11): 2717–2734, doi:10.1007/s11229-019-02173-1.
- 2215 GAFFNEY, Jeffrey and MARLEY, Nancy. 2018. *General Chemistry for Engineers*. Ams-  
 2216 terdam: Elsevier Science Publishers B.V., doi:10.1016/C2015-0-05956-1.
- 2217 GAO, Shan. 2017. *The Meaning of the Wave Function: In Search of the Ontology of*  
 2218 *Quantum Mechanics*. Cambridge: Cambridge University Press, doi:10.1017/9781  
 2219 316407479.
- 2220 HEIL, John. 2003. *From an Ontological Point of View*. Oxford: Oxford University Press,  
 2221 doi:10.1093/0199259747.001.0001.
- 2222 —. 2012. *The Universe as We Find It*. Oxford: Oxford University Press, doi:10.1093/ac  
 2223 prof:oso/9780199596201.001.0001.
- 2224 JANSEN, Charles M. 2019. "De-Fining Material Things." *Dialectica* 73(4): 459–477,  
 2225 doi:10.1111/1746-8361.12280.
- 2226 JAWORSKI, William. 2016. *Structure and the Metaphysics of Mind: How Hylomorphism*  
 2227 *Solves the Mind-Body Problem*. Oxford: Oxford University Press, doi:10.1093/acprof:  
 2228 oso/9780198749561.001.0001.
- 2229 JENKINS, Frank and VAN KESSEL, Hans. 2002. *Nelson Chemistry 12*. Toronto: Thomson  
 2230 Nelson.
- 2231 KEINÄNEN, Markku. 2018. "Instantiation and Characterization: Problems in Lowe's  
 2232 Four-Category Ontology." in *Studies in the Ontology of E.J. Lowe*, edited by Timothy  
 2233 TAMBASSI, pp. 109–124. Neunkirchen-Seelscheid: editiones scholasticae.
- 2234 KEINÄNEN, Markku and HAKKARAINEN, Jani. 2017. "Kind Instantiation and Kind  
 2235 Change: A Problem for Four-Category Ontology." *Studia Neoaristotelica* 14(2):  
 2236 139–165, doi:10.5840/studneoar20171427.
- 2237 KOONS, Robert C. 2018. "Hylomorphic Escalation: An Aristotelian Interpretation  
 2238 of Quantum Thermodynamics and Chemistry." *American Catholic Philosophical*  
 2239 *Quarterly* 92(1): 159–178, doi:10.5840/acpq2017124139.
- 2240 KOSLICKI, Kathrin. 2008. *The Structure of Objects*. Oxford: Oxford University Press,  
 2241 doi:10.1093/acprof:oso/9780199539895.001.0001.
- 2242 —. 2018. *Form, Matter, Substance*. Oxford: Oxford University Press, doi:10.1093/oso/  
 2243 9780198823803.001.0001.
- 2244 LOUX, Michael J. 1998. *Metaphysics: A Contemporary Introduction*. Routledge Con-  
 2245 temporary Introductions to Philosophy. London: Routledge, doi:10.4324/97802034  
 2246 38244.
- 2247 —. 2006. *Metaphysics: A Contemporary Introduction*. 3rd ed. Routledge Contemporary  
 2248 Introductions to Philosophy. London: Routledge. Revised edition of Loux (1998),  
 2249 doi:10.4324/9780203968871.
- 2250 MARKOSIAN, Ned. 2008. "Restricted Composition." in *Contemporary Debates in Meta-*  
 2251 *physics*, edited by Theodore SIDER, John HAWTHORNE, and Dean W. ZIMMERMAN,  
 2252 pp. 341–364. Contemporary Debates in Philosophy n. 10. Malden, Massachusetts:  
 2253 Wiley-Blackwell.

- 2254 MARMODORO, Anna and PAGE, Ben. 2016. "Aquinas on Forms, Substances and Arti-  
2255 facts." *Vivarium* 54(1): 1–21, doi:10.1163/15685349-12341310.
- 2256 MCKENZIE, Kerry. 2011. "Arguing Against Fundamentality." *Studies in History and*  
2257 *Philosophy of Science. Part B: Studies in History and Philosophy of Modern Physics*  
2258 42(4): 244–255, doi:10.1016/j.shpsb.2011.09.002.
- 2259 —. 2014. "Priority and Particle Physics: Ontic Structural Realism as a Fundamentality  
2260 Thesis." *The British Journal for the Philosophy of Science* 65(2): 353–380, doi:10.109  
2261 3/bjps/axt017.
- 2262 MCKENZIE, Kerry and MULLER, F. A. 2017. "Bound States and the Special Composition  
2263 Question." in *EPSA15 Selected Papers: The 5th Conference of the European Philoso-*  
2264 *phy of Science Association in Düsseldorf*, edited by Michela MASSIMI, Jan-Willem  
2265 ROMELIN, and Gerhard SCHURZ, pp. 233–241. European Studies in Philosophy of  
2266 Science n. 5. Cham: Springer, doi:10.1007/978-3-319-53730-6\_19.
- 2267 NEY, Alyssa and ALBERT, David Z., eds. 2013. *The Wave Function: Essays on the*  
2268 *Metaphysics of Quantum Mechanics*. Oxford: Oxford University Press, doi:10.1093/  
2269 aprof:oso/9780199790807.001.0001.
- 2270 PASNAU, Robert. 2011. *Metaphysical Themes 1274–1671*. Oxford: Oxford University  
2271 Press, doi:10.1093/acprof:oso/9780199567911.001.0001.
- 2272 —. 2012. "Mind and Hylomorphism." in *The Oxford Handbook of Medieval Philosophy*,  
2273 edited by John MARENBO, pp. 486–504. Oxford Handbooks. Oxford: Oxford  
2274 University Press, doi:10.1093/oxfordhb/9780195379488.013.0022.
- 2275 PAWL, Timothy and SPENCER, Mark K. 2016. "Christologically Inspired, Empirically  
2276 Motivated Hylomorphism." *Res Philosophica* 93(1): 137–160, doi:10.11612/resphil  
2277 .2016.93.1.6.
- 2278 RENZ, Graham. 2018. "Form as Structure: It's Not So Simple." *Ratio* 31(1): 20–36,  
2279 doi:10.1111/rati.12155.
- 2280 ROONEY, James Dominic. 2022. *Material Objects in Confucian and Aristotelian Meta-*  
2281 *physics: The Inevitability of Hylomorphism*. London: Bloomsbury Academic, doi:10  
2282 .5040/9781350276376.
- 2283 SCALTSAS, Theodore. 1994. *Substances and Universals in Aristotle's Metaphysics*. Ithaca,  
2284 New York: Cornell University Press.
- 2285 STUMP, Eleonore. 2003. *Aquinas. The Arguments of the Philosophers*. London: Rout-  
2286 ledge, doi:10.4324/9780203928356.
- 2287 THOMAS AQUINAS. 1912–1936. *Summa Theologica*. London: Burns, Oates & Wash-  
2288 bourne. Translated by the Fathers of the English Dominican Province, <http://www.domcentral.org/summa/>.
- 2289 —. 1926. *De ente et essentia*. Opuscula et textus. Series scholastica n. 1. Münster i.W.:  
2290 Monasterii; Aschendorf. Edidit Ludovicus Baur.
- 2291 —. 1949. *The Soul: A Translation of St. Thomas Aquinas' De anima*. St. Louis: B. Herder  
2292 Book Co. Transl. by John Patrick Rowan, ark:/13960/t87j03ros.
- 2293

- 2294 —. 1950. *Sententia libri Metaphysicae*. Roma: Domus Editorialis Marietti Taurini.  
 2295 Translated and edited by Roberto Busa, <https://www.corpusthomicum.org/cm>  
 2296 [p00.html](https://www.corpusthomicum.org/cm/p00.html).
- 2297 —. 1953. *Quaestiones disputatae de anima*. Roma: Domus Editorialis Marietti Taurini.  
 2298 Translated and edited by Roberto Busa, <https://www.corpusthomicum.org/qd>  
 2299 [a01.html](https://www.corpusthomicum.org/qd/a01.html).
- 2300 —. 1956. *Quaestiones Quodlibetales*. Roma: Domus Editorialis Marietti Taurini. Trans-  
 2301 lated and edited by Raimondo Spiazzi, <https://www.corpusthomicum.org/qda0>  
 2302 [1.html](https://www.corpusthomicum.org/qda0/1.html).
- 2303 —. 1961a. *Commentary on the Metaphysics of Aristotle, Volume 1*. Library of Living  
 2304 Catholic Thought. Chicago, Illinois: Henry Regnery Company. Transl. by John  
 2305 Patrick Rowan, [ark:/13960/t5gb5411f](https://www.corpusthomicum.org/ark:/13960/t5gb5411f).
- 2306 —. 1961b. *Commentary on the Metaphysics of Aristotle, Volume 2 (Books VI-XII)*. Library  
 2307 of Living Catholic Thought. Chicago, Illinois: Henry Regnery Company. Transl. by  
 2308 John Patrick Rowan, [ark:/13960/t5gb5411f](https://www.corpusthomicum.org/ark:/13960/t5gb5411f).
- 2309 —. 1968. “De substantiis separatis.” in, pp. D1–D87. Leonina, XL. Roma: Ad Sanctae  
 2310 Sabinae. Ed. H.-F. Dondaine, <https://www.corpusthomicum.org/ots.html>.
- 2311 —. 1976. “De principiis naturae [ad Fratrem Sylvestrum].” in *Opuscula IV*, pp. 1–47.  
 2312 Leonina, XLIII. Roma: Editori di San Tommaso. Ed. H.-F. Dondaine, based on the  
 2313 previous work of P.J.M. Perrier.
- 2314 —. 1983. *Quodlibetal Questions 1 and 2*. Mediaeval Sources in Translation n. 27. Toronto:  
 2315 Pontifical Institute of Mediaeval Studies. Translated with an introduction and  
 2316 notes by Sandra Edwards.
- 2317 —. 1996a. *Quaestiones disputatae de anima*. Leonina, XXIV, 1. Paris: Éditions du Cerf.  
 2318 Ed. B.C. Bazán.
- 2319 —. 1996b. *Quaestiones de quolibet. Préface, Quodl. VII, VIII, IX, X, XI*. Leonina, XXV, 1.  
 2320 Paris: Éditions du Cerf. Ed. R.-A. Gauthier.
- 2321 —. 1996c. *Quaestiones de quolibet. Quodl. I, II, III, VI, IV, V, XII*. Leonina, XXV, 2. Paris:  
 2322 Éditions du Cerf. Ed. R.-A. Gauthier.
- 2323 —. n.d. *Summa Theologica*. Leonina, IV-XII. Roma: Typographia poliglotta S.C. de  
 2324 Propaganda Fide.
- 2325 TONER, Patrick. 2008. “Emergent Substance.” *Philosophical Studies* 141(3): 281–297,  
 2326 doi:[10.1007/s11098-007-9160-6](https://doi.org/10.1007/s11098-007-9160-6).
- 2327 WIPPEL, John F. 2000. *The Metaphysical Thought of Thomas Aquinas: From Finite*  
 2328 *Being to Uncreated Being*. Washington, D.C.: The Catholic University of America  
 2329 Press.
- 2330 —. 2011. “Thomas Aquinas and the Unity of Substantial Form.” in *Philosophy and*  
 2331 *Theology in the Long Middle Ages: A Tribute to Stephen F. Brown*, edited by Kent  
 2332 EMERY Jr., Russell L. FRIEDMAN, and Andreas SPEER, pp. 117–154. Studien und  
 2333 Texte zur Geistesgeschichte des Mittelalters n. 105. Leiden: E.J. Brill, doi:[10.1163/](https://doi.org/10.1163/ej.9789004169425.i-1006.32)  
 2334 [ej.9789004169425.i-1006.32](https://doi.org/10.1163/ej.9789004169425.i-1006.32).

# A Generalization of the Reflection Principle

WOLFGANG SPOHN

This paper generalizes (probabilistic) auto-epistemology by amending the original forward-looking reflection principle of van Fraassen (1984), which is about learning or favorable epistemic changes in general, by a new but similar backward-looking principle, which is about forgetting and other unfavorable epistemic changes. The generalization is argued to be a completion by defending what is called the no-neutrality condition. Due to the similarity, analogous consequences are provable for both principles. This fact is utilized for a plausibility check of the new principle. Finally, it is argued that this generalization should not be considered as a special case of the expert principle.

Van Fraassen (1984, 244) has introduced the *reflection principle*, as he called it: *Given your tomorrow's probabilities are such and such, your present conditional probabilities should be the very same.* Hence, this principle may be called *forward-looking*. It is *the* basic principle of probabilistic or Bayesian auto-epistemology. Various interesting consequences have been derived from it. It has met diverse criticisms, several variants have been offered in response, and it is well-known by now that it holds only under restrictions. However, the literature does not offer clear ideas on whether and how it may be suitably amended or even completed. This paper tries to do better by proposing what I call the full reflection principle. It has certainly been suggested in one way or another. But I am after precise statements allowing strict inferences.

Here is how the paper proceeds: Section 1 briefly recapitulates van Fraassen's original principle. Section 2 suggests an equally strong and formally analogous *backward-looking* principle. Both combine to what I call the full reflection principle. Thereby, I propose to double, as it were, and arguably complete the range of auto-epistemology. In section 3, the completeness claim is supported by what I call the no-neutrality condition. However, the completeness claim does not go so far as to offer an account

2365 of how to iteratively apply the full principle. Section 4 looks at some  
 2366 well-known consequences of van Fraassen's principle, which carry over to  
 2367 the backward-looking principle. This is intended to serve as a plausibility  
 2368 check of the proposed generalization. Section 5 discusses what happens  
 2369 if the reflection principles are referred to as auto-epistemic propositions  
 2370 themselves; this is often not strictly distinguished. There, we will discover a  
 2371 slight disanalogy between the forward- and the backward-looking principle.  
 2372 Section 6, finally, defends my generalization against the objection that it is  
 2373 already contained in the familiar generalization of van Fraassen's principle,  
 2374 known as the expert principle.

### 2375 **1 Van Fraassen's Reflection Principle**

2376 This paper intends to focus on the philosophical and not on the formal aspects  
 2377 of its topic. Still, the main claims should be stated in a formally precise way.  
 2378 For this purpose, we need a modicum of notation. We will refer to a prior time,  
 2379 0, and a posterior time, 1. This may be today and tomorrow, or yesterday and  
 2380 today; we will use both readings. Everything indexed in this way will refer so  
 2381 as well. Thus,  $P_0$  is to represent your actual prior and  $P_1$  your actual posterior  
 2382 credences, where I suppose you to satisfy all rationality constraints we think to  
 2383 be pertinent. These are at least the synchronic axioms of probability and some  
 2384 diachronic learning rules we need not fix. Without an index, the temporal  
 2385 reference may be any. I use  $\pi$ , with or without indices, as a variable for your  
 2386 possible probability functions. Sets of those functions represent *auto-epistemic*  
 2387 propositions about your own probabilities; e.g.,  $\{\pi_0 \mid \pi_0(A) = x\}$  represents  
 2388 the proposition that your prior probability for  $A$  is  $x$ ,  $\{\pi_1 \mid \pi_1 = Q_1\}$  represents  
 2389 the proposition that your posterior probability function is  $Q_1$ , etc.

2390 These possible probability functions are about some fixed algebra of propo-  
 2391 sitions concerning some worldly affairs. Not all worldly affairs; the range may  
 2392 be quite restricted but need not be made explicit. The label "worldly" is to  
 2393 mean that these propositions are about external matters and not about your  
 2394 epistemic states. (In section 5, we will consider dropping this restriction.) It  
 2395 also means that the propositions are about empirical, not abstract or formal,  
 2396 matters. The epistemology of formal sciences is a very different topic not  
 2397 amenable to the present methods in my view.

2398 However, your actual  $P_0$  and  $P_1$  are not only about these worldly propositions  
 2399 but also about all the auto-epistemic propositions just introduced and all

2400 algebraic combinations thereof. Thus,  $P_0$  and  $P_1$ , but not the possible  $\pi$ , reflect  
 2401 on how your prior and posterior probabilities might be.

2402 Now, we are ready to explicitly state van Fraassen's principle; I will call it  
 2403 the *forward-looking reflection principle* for quite obvious reasons:

2404 (1)  $P_0(A \mid \{\pi_1 \mid \pi_1 = P_1\}) = P_1(A)$  for all worldly propositions  $A$ .

2405 That is, *given* your posterior credence is  $P_1$  and, in particular,  $P_1(A)$  for  $A$ ,  
 2406 your prior for  $A$  is also  $P_1(A)$ .<sup>1</sup> Why do I state the condition in (1) as  $\{\pi_1 \mid$   
 2407  $\pi_1 = P_1\}$  and not in the usual way as  $\{\pi_1 \mid \pi_1(A) = P_1(A)\}$ ? Because the usual  
 2408 version is weaker and derivable from the present version, which is clearly  
 2409 the intended one.<sup>2</sup> Note, moreover, that " $\pi_1 = P_1$ " cannot be literally true,  
 2410 given that  $\pi_1$  is only about worldly propositions, while  $P_1$  is also about auto-  
 2411 epistemic propositions. Strictly speaking, I should refer to  $P_1$  as restricted  
 2412 to worldly propositions. This would be too cumbersome, though. This little  
 2413 sloppiness will do no harm.

2414 Let me emphasize that I am using probabilistic terminology here only for  
 2415 convenience; it is the most familiar one. However, the reflection principle is  
 2416 not restricted to Bayesian epistemology. It holds for any kind of epistemic for-  
 2417 mat that allows for conditional epistemic states. Binkley (1968) has proposed  
 2418 a qualitative version: if you now believe that you will believe  $p$  tomorrow,  
 2419 you should believe  $p$  already now. This principle played a crucial role in his  
 2420 account of the surprise examination paradox. In Spohn (2010, 125–127), I  
 2421 have stated the reflection principle in terms of ranking theory. One may state  
 2422 the reflection principle in terms of imprecise or interval probabilities. And  
 2423 so on.<sup>3</sup> Here, I will simply assume that epistemic states are represented in a  
 2424 fixed formal format, and then I choose the most familiar one. I do not want

1 In order to keep things simple, this statement is predicated on the assumption that there are only finitely many possible  $\pi$  under consideration. Generally, though, there are uncountably many  $\pi$  in play, and then the condition is likely to have probability 0 so that the conditional probability in (1) is undefined. Then we must replace (1) by Constraint 2, as Skyrms (1980, 163) calls it, which says for all intervals  $I$ :  $P_0(A \mid \{\pi_1 \mid \pi_1(A) \in I\}) \in I$ . This solves the problem because now the condition will usually have a positive probability. See also Goldstein (1983).

2 For instance, from (1), but not from the usual version, one can derive  $P_0(A \mid \{\pi_1 \mid \pi_1(A) = x \ \& \ \pi_1(B) = y\}) = x$ , which is desirable, too.

3 Schoenfeld (2012) argues that the reflection principle *entails* probabilities to be precise—and is, therefore, to be rejected because other arguments require probabilities to be imprecise. If so, one should state, I think, a reflection principle for interval probabilities:  $P_0(A \mid \{\pi_1 \mid \pi_1(A) = I\}) = I$ , where  $P_0$  and the  $\pi_1$  would now be interval probability functions. (Observe the difference to "Constraint 2" in footnote 1.)

to discuss the more basic issue of the most adequate representation of epistemic states. Hence, the entire paper will move within a standard probabilistic representation of epistemic states, though only *pars pro toto*.

I shall not repeat here the grounds on which (1) is widely accepted; I shall only refer to some criticism below. I should mention, though, that (1) has been considered before. Spohn (1978, 161–162) and Goldstein (1983) have proposed the iteration principle (which is almost equivalent; see sections 4 and 5). Quite generally, one might say that the reflection principle is implicit or well-nigh explicit in de Finetti's (1937) philosophy of probability. One incurs considerable philosophical costs if one wants to abandon it.

Gaifman (1986) has suggested a more general reading of (1). He defines  $P_1$  to be an expert for you at the prior time 0 (regarding  $A$ ) if and only if (1) holds for your  $P_0$ . In this reading, (1) turns into what is called the *expert principle*, saying that you should trust the experts—though this is now tautological because an expert has been defined as one you can trust in this way. Still, the maneuver provides quite a graphic reading of the original (1), namely, as saying that you should accept your posterior opinion as an expert for you.<sup>4</sup> In section 6, I shall discuss whether the amendment of (1) proposed in section 2 may simply be conceived as an instance of this expert principle.

This makes (1) plausible. Surely, if your posterior probabilities have learned something, have gathered evidence, are better informed, etc., then your prior probabilities should consider them to be an expert for you. At the same time, this points to the restricted applicability of (1). Your future self is not always better informed. You may forget things, you may be tired, brainwashed, confused by drugs, your judgment may be obfuscated by prejudices against better knowledge, etc. In all those cases, your future opinion is not a trustworthy expert for you. This restriction has often been noticed—e.g., by Christensen (1991) (referring to epistemic change due to psychedelic drugs), by Talbott (1991) (referring mainly to memory loss), or by Spohn (1978, 166) (with re-

---

4 The striking similarity between the reflection principle and Lewis's Principal Principle has often been noticed; see, e.g., Spohn (2010). In these terms, the Principal Principle says that chance at  $t$  is your best expert at  $t$ , and truth = chance at the end of time is your best expert at all. Christensen (2010b) discusses a principle of rational reflection, which is about your probabilities being conditional not on your future probabilities but on your current probabilities as they should rationally be. This differs from the original reflection principle only in case you suspect your future probabilities to diverge from how the current probabilities should rationally be. In any case, I stick to the original principle and won't discuss such subtly related principles. The real challenge raised by Gaifman is, of course, what to do when I am impressed by several (diverging) experts. Again, peer disagreement is not my topic here.



2454 spect to the iteration principle). Briggs (2009, 64ff.) presents a nice list of  
2455 exceptions.<sup>5</sup>

2456 In my view, this restriction does not diminish the importance of the reflec-  
2457 tion principle (1). Whenever your epistemic changes, be they rational or not,  
2458 do not just occur to you and you rather take a reflective attitude towards them,  
2459 (1) serves as a fundamental meta-principle. It does not specify how learning  
2460 precisely works. There, we may consider rules like simple or Jeffrey condi-  
2461 tionalization, minimizing relative entropy, etc. But it says, however learning  
2462 works, it must satisfy (1) if reflected.

2463 Surely, though, (simple) conditionalization is the basic Bayesian learning  
2464 rule. Since van Fraassen (1984) proposed a Dutch book argument in favor of  
2465 his reflection principle, critical discussions like those of Christensen (1991)  
2466 and Talbott (1991) focused on the tension of this argument with the familiar  
2467 Dutch book justification of conditionalization. Indeed, the relation between  
2468 the reflection principle and the conditionalization rule is delicate, as displayed  
2469 also in Weisberg (2007). However, Hild (1998a, 1998b) has already shown  
2470 that the reflection principle is equivalent to a rule he calls auto-epistemic  
2471 conditionalization, and he has fully stated the conditions under which auto-  
2472 epistemic conditionalization and simple and Jeffrey conditionalization may  
2473 come apart. Therefore, I shall not pursue this connection any further.

## 2472 **A Generalization of Van Fraassen's Reflection Principle**

2475 If the forward-looking principle (1) has only restricted validity, we should  
2476 think about whether there are auto-epistemic principles governing the cases  
2477 not covered by (1). These may be irrational epistemic changes violating di-  
2478 achronic principles of rationality that are supposed to be governed by the  
2479 reflection principle (1), whatever they are. E.g., whenever your prejudices  
2480 get a hold on you, although you are sure that they are irrelevant, this is ar-  
2481 guably a case of epistemic irrationality. However, there may also be arational  
2482 epistemic changes, which do not violate diachronic principles but simply  
2483 fall outside the scope of such principles, such as fatigue or clouding one's  
2484 judgment by getting drunk (where the only irrationality may have been to get  
2485 drunk), etc. Forgetting is perhaps the most familiar case of such an arational  
2486 change. However, it consists not only in unlearning a fact. It can take many

---

5 Van Fraassen (1995) takes the strange recourse to say that in such cases, your epistemic self ceases to exist, thus recovering the general validity of (1) for you.

2487 shapes. E.g., one may forget live possibilities so that one becomes sure of the  
 2488 remaining possibilities. Or one can forget about evidential relations, say, that  
 2489 a characteristic smell is a sign of a poison. And so on.

2490 Such changes seem to be a matter of empirical psychology and not of a  
 2491 normative theory of epistemic rationality. Isn't it an empirical question what  
 2492 we are prone to forget, how our beliefs are influenced by prejudices or drugs,  
 2493 etc.? Yes. However, this does not mean that a theory of epistemic rationality  
 2494 cannot say anything about how we should *rationally deal* with such arational  
 2495 and irrational changes when we, possibly falsely, suspect them to occur. On  
 2496 the contrary, if we are thinking about how to extend the reflection principle  
 2497 (1), this is the challenge we positively face. Hence, the question is: Is there  
 2498 an auto-epistemic rationality principle dealing also with those irrational and  
 2499 arational changes?

2500 Yes, there is. I follow the idea of Titelbaum (2013, chap. 6): just reverse the  
 2501 temporal perspective!<sup>6</sup> Place yourself at your posterior  $P_1$  and ask whether you  
 2502 should consult your prior  $P_0$ , whatever it was. Certainly not in van Fraassen's  
 2503 cases where you are better informed in  $P_1$ . But surely in those cases where you  
 2504 have forgotten something, are foggy-brained, etc. in  $P_1$ . Therefore, I propose  
 2505 the following *backward-looking reflection principle*:<sup>7</sup>

2506 (2)  $P_1(A \mid \{\pi_0 \mid \pi_0 = P_0\}) = P_0(A)$  for all worldly propositions  $A$ .<sup>8</sup>

6 I have first proposed an essential part of the following considerations in Spohn (2017, sec. 10) in a more convoluted way at the end of a long investigation on indexical belief. The connection was Sleeping Beauty, which seems to be not only a problem about indexical belief but also about auto-epistemology. It has been observed, e.g., by Arntzenius (2003) or Mahtani (2017), that the proposition on which the probabilities in the reflection principle are conditioned must egocentrically refer to *me* and *my future* probabilities, not to *a*'s probabilities at *t*, where *a* happens to be me and *t* happens to lie in the future. Here, I ignore this line of thought. The subject's potential uncertainty about her own location is not our issue.

7 Christensen (2000, 352ff.) also speaks of taking a backward-looking perspective, but he thereby means something else. He is interested in diachronic coherence in the sense of an epistemic conservatism, which seeks to preserve as much of the old beliefs as possible while *learning* something new.

8 As just mentioned, Titelbaum (2013) has already observed the symmetry between the forward- and the backward-looking case. He captured both in his principle of Generalized Conditionalization (Titelbaum 2013, 127). I should have known and noticed this in Spohn (2017). However, he refers his observation to simple conditionalization and conceives of the prior state as the conditionalization of the posterior forgetful state with respect to the forgotten proposition. On p. 133, he also considers a generalized reflection principle, but only as derived from simple conditionalization. By contrast, I think of the reflection principles as having an independent role and apply the backward-looking principle to all kinds of arational and irrational epistemic

2507 (2) is as important and fundamental a meta-principle as (1). Whatever the  
 2508 multifarious changes of our epistemic position to the worse, if we reflect on  
 2509 them, (2) must be obeyed in all such cases. In a decision-theoretic perspective,  
 2510 which I do not unfold in this paper, this meta-principle is at the bottom of our  
 2511 efforts to fight forgetting, e.g., by building museums and archives and even  
 2512 inventing scripture, and of our attempts to preserve our epistemic integrity  
 2513 wherever we can, e.g., by banning brainwashing.

2514 It is clear, however, that not both (1) and (2) can be applied generally.  
 2515 Their applications cannot even overlap. In a case of such overlap, your prior  
 2516 probability should trust your posterior one, but your posterior one should  
 2517 reversely follow your prior one; so, if both are mutually envisaged, they must  
 2518 be the same. In other words, only in the case of non-change can (1) and (2)  
 2519 apply simultaneously. In section 5, I will provide a formal proof of this claim.

2520 Hence, the fields of application of (1) and (2) must be disjoint (with the  
 2521 possible exception of non-change). Indeed, this is how I have explained (1)  
 2522 and (2). In order to have uniform labels, let's say that (1) applies in the case of  
 2523 *favorable* changes, such as learning, receiving information, acquiring evidence,  
 2524 etc., and that (2) applies in the case of *unfavorable* changes, such as forgetting,  
 2525 drinking too much, being influenced by prejudices one takes to be unjustified,  
 2526 etc. In a case of a favorable change from  $P$  to  $P'$  or an unfavorable change  
 2527 from  $P'$  to  $P$ , let's say that  $P'$  is *superior* to  $P$  and  $P$  is *inferior* to  $P'$ .

2528 The next question is: Who is to judge changes as favorable or unfavorable  
 2529 and epistemic positions as superior or inferior? We who decree these principles  
 2530 from outside? No, I think it is preferable to subjectivize the application  
 2531 conditions of (1) and (2). The subject herself must assess the epistemic changes  
 2532 she is considering: Is my change from  $P_0$  to  $P_1$  favorable and  $P_1$  superior? Then  
 2533 apply (1)! Is my change from  $P_0$  to  $P_1$  unfavorable and  $P_0$  superior? Then apply  
 2534 (2)! The first instance to assess this is the subject herself.

2535 Of course, this does not preclude that, in a second step, we have a normative  
 2536 argument about this assessment. Presumably, we all agree that experience  
 2537 is favorable and forgetting is unfavorable. But what about hunches and gut  
 2538 feelings? Gigerenzer (2007) is a strong plea for respecting gut feelings not only  
 2539 as a psychological fact but also as a guideline to rational decision-making.  
 2540 Prejudices may not always be bad. What about epiphanies? Those claiming  
 2541 to have had them feel to be in a superior position. Surely, these are difficult

---

changes, not just to the forgetting of previous certainties, as Titelbaum does by only dealing with simple conditionalization. Insofar my approach is more general.

2542 and possibly contested issues that we may and must discuss. However, the  
 2543 reflection principles as such are independent of that discussion, and, therefore,  
 2544 we should keep matters separate. The intent of my subjectivizing move was  
 2545 precisely not to get involved in that discussion.

2546 Let's take a slightly more general perspective for stating it. The temporal  
 2547 relations do not really seem to matter. The point is rather that in whatever  
 2548 epistemic position I am, I would trust a superior position and mistrust an  
 2549 inferior position (where it is up to me what I take to be superior and inferior).  
 2550 This seems to be the gist of the principle. If so, we arrive at the following  
 2551 *full reflection principle*, in which the temporal location of the probabilities  
 2552 referred to is left open (hence no indices):

2553 (3) for all worldly propositions  $A$   $P(A \mid \{\pi \mid \pi = P'\}) = Q(A)$ , given that  $Q$   
 2554 is taken to be the superior one of  $P$  and  $P'$ .<sup>9</sup>

2555 This is the generalization of van Fraassen's auto-epistemology I would like to  
 2556 propose.

### 2553 3 The No-Neutrality Condition and the Iteration Problem

2558 Is the generalization a completion? This raises two issues. First, the principle  
 2559 (3) reflects only upon a single possible change. But we may certainly reflect  
 2560 on iterated change. (3) is silent on this and insofar still incomplete.<sup>10</sup> I will  
 2561 not be able to offer a solution, but I will comment on the issue below.

2562 Secondly, if we attend only to a single change, the reflection principle (3)  
 2563 is complete only if the *no-neutrality condition* holds, which states that there  
 2564 are no neutral and no incomparable changes; there are only either favorable  
 2565 or unfavorable changes and nothing else. Then, but only then, there would

---

9 The "is taken" always refers to the assessment of the subject we are talking about, not our own. In Spohn (2017), I have emphasized this by making the condition in (3) part of the conditional probability statement. However, this raises awkward questions. Are propositions of the form " $P'$  is superior to  $P$ " part of the auto-epistemically extended algebra of propositions? Do they receive probabilities? Are these probabilities subject to change? We better avoid such questions. These propositions are outside the epistemic game we are considering. We may rather assume that the subject's superiority assessments are stable within our dynamic scenario. Therefore, I now state this condition outside the probability statement, though still in a subjectivized form.

10 The importance of this issue is underscored by the parallel case in belief revision theory, which was initially restricted to treating only single revisions and thus plagued by the iteration problem, too. I have first raised it in Spohn (1988, 112ff.). It turned out to virtually be an anomaly in the Kuhnian sense; see, e.g., Rott (2009).

2566 always be *the superior one* of  $P$  and  $P'$  (with the irrelevant exception of non-  
2567 change), and the application condition of (3) would be complete. Does this  
2568 condition hold?

2569 Yes, I think so. I welcome favorable changes and seek superior epistemic  
2570 positions (if they are not too costly), and I try to prevent unfavorable changes  
2571 and avoid inferior positions (if that is not too costly, either). At least this is  
2572 so by purely epistemic standards; moral standards, e.g., may tell otherwise in  
2573 special cases. Thus, a change that is neither favorable nor unfavorable would  
2574 be one I don't care about. I would say then: it's nice to have the present prior  
2575  $P_0$ , and it's equally nice to have the posterior  $P_1 \neq P_0$  later on; both are fine,  
2576 and none of them is inferior or superior. This sounds very strange to me. This  
2577 makes the change from  $P_0$  to  $P_1$  appear arbitrary and without good reason, and  
2578 then I can't stay indifferent about the change; it must appear unfavorable to  
2579 me.

2580 To illustrate: Today, I think I will be in good health next year, and tomorrow,  
2581 just after sleeping and without any new information whatsoever, I think  
2582 I won't. Usually, this would not be taken as a change of mind but as an  
2583 expression of continued uncertainty. But say, today, I am firmly convinced  
2584 that I will be in good health next year. From this perspective, it must appear  
2585 arbitrary when I would have changed my mind by tomorrow to equally firmly  
2586 believe the contrary. It would be odd to presently be neutral about such a  
2587 change; I should rather reject and not trust it.<sup>11</sup>

2588 This is not a cogent argument. It is only to say that I cannot imagine how  
2589 the no-neutrality condition could be violated. In any case, one must be aware  
2590 that this condition is a crucial and substantial normative principle. If we  
2591 accept it, then (3) indeed deserves the label "full reflection principle," at least  
2592 regarding single changes.<sup>12</sup>

2593 I think, though, that there is a deeper reason behind the no-neutrality  
2594 condition. It is that, ultimately, there is only one standard for our epistemic  
2595 states: truth. We try to approach the truth and avoid veering away from the  
2596 truth, however we measure the distance here. The point is that there is only  
2597 one "scale" to measure. If epistemic states had to meet many standards on  
2598 different scales, then indifferences or even incomparabilities might easily  
2599 arise. Such more complex situations would certainly be relevant when we  
2600 were to more generally think about what kind of person we want to be. There

---

11 See also the arguments against arbitrary switching in White (2014, 318ff.).

12 As mentioned, the case of non-change from  $P$  to  $P'$  may be ruled arbitrarily. We may say then that  $P$  is superior, or  $P'$  is, or both are. It doesn't make any difference for (3).

2601 are many aspects in which we change for the better or the worse, and we  
 2602 will often have indeterminate preferences about possible personal changes or  
 2603 none at all. But in the case of epistemic change, our judgments seem to be  
 2604 unambiguous.

2605 So far, I have only argued that there are no neutral changes. And I have  
 2606 excluded the possibility of incomparabilities due to a multitude of epistemic  
 2607 standards. However, there are easier ways for incomparabilities to arise. Surely,  
 2608 there are complex changes that are favorable in some respects and unfavor-  
 2609 able in others so that, overall, the result is neither superior nor inferior but  
 2610 incomparable. For instance, I learn that I have a date with the president next  
 2611 Friday and simultaneously forget that I have already agreed to meet the vice  
 2612 president at the very same time. I propose to treat this as two changes, first a  
 2613 favorable one and then an unfavorable one—or the other way around; it is not  
 2614 guaranteed that this comes to the same. Often, a temporal succession can be  
 2615 discerned within such a complex change, and sometimes, e.g., in my example,  
 2616 this move may be artificial. However, my proposal seems feasible; it avoids the  
 2617 need to refer favorability and unfavorability to aspects of complex changes,  
 2618 and it saves the no-neutrality condition. So, in any case, it is theoretically  
 2619 beneficial.

2620 However, this move makes the first issue of extending the full reflection  
 2621 principle (3) to iterated change more pressing. To my knowledge, this issue  
 2622 has not been considered in the literature. Perhaps the reason is that it seemed  
 2623 trivial in the case of the original reflection principle. If my first epistemic  
 2624 state trusts the second, and the second trusts the third, already the first state  
 2625 can trust the third. Reversely, in the case of iterated forgetting.<sup>13</sup> However,  
 2626 there are also mixed cases, and I have just alluded to them.<sup>14</sup> The difficult  
 2627 case is the one where my epistemic state first changes in an unfavorable way  
 2628 and then in a favorable way; e.g., first, I forget some things, and then I learn  
 2629 other things. In this case, my initial state can neither trust the final state in  
 2630 the sense of principle (1) nor dismiss the final state in the sense of principle  
 2631 (2). Rather, it seems that I have to engage in a counterfactual consideration. In  
 2632 this case, I can only trust that the epistemic state that *would have* emerged had  
 2633 I not incurred the first unfavorable change (forgotten the one thing) but still

13 Titelbaum (2013, 140) seems to be able to treat the iterated case with the help of his principle of suppositional consistency. But if so, this is due to the fact that the only epistemic changes he considers are the gain and loss of certainties.

14 I have discussed the various cases and their problems a bit more extensively in Spohn (2017, 408–409).

2634 experienced the second favorable change (learned the other things). That is, I  
 2635 would have to speculate not only about my actual epistemic states and their  
 2636 change but also about my counterfactual epistemic states and their change.  
 2637 Hence, a general solution of this problem seems to require quite different  
 2638 theoretical means. It is not a task we can pursue here.

2639 Still, it should be pursued. To emphasize its urgency: As far as I see, the issue  
 2640 of so-called second- or higher-order evidence is closely related. Christensen  
 2641 (2010a) gives a wide variety of examples. A salient structure of them (not all  
 2642 of them) is this: I receive a lot of ordinary (first-order) evidence on a certain  
 2643 matter, and I seem to draw my conclusions from it in the usual rational way.  
 2644 At the same time, I receive higher-order evidence (perhaps falsely) indicating  
 2645 that my cognitive abilities are somehow hampered. I am overly tired, I am  
 2646 told to have consumed a fancy drug, I am instructed that I regularly tend  
 2647 to overoptimism, I may be suffering from hypoxia (a realistic example from  
 2648 Christensen 2010b, 126), etc. So, maybe I should correct my inferences?

2649 In such cases, the higher-order evidence indicates that I should not trust  
 2650 the epistemic state I have reached. But, neither can I simply rely on my  
 2651 prior epistemic state before the change, as the backward-looking reflection  
 2652 principle (2) would have it. As mentioned above, such cases are mixtures  
 2653 of two different epistemic movements. On the one hand, there is the first-  
 2654 order evidence, which I should trust. On the other hand, there is my alleged  
 2655 epistemic handicap, which suggests an accompanying unfavorable change and  
 2656 should make me think about what my inferences would have been without  
 2657 the handicap. So, the issue of higher-order evidence would also profit from  
 2658 solving the iteration problem. However, for the reasons indicated, I don't  
 2659 address this problem further. We should be content with treating the pure  
 2660 cases.<sup>15</sup>

#### 2664 **Are the Consequences of the Full Principle Acceptable?**

2662 The full reflection principle (3) seems to be intuitively plausible and philo-  
 2663 sophically important. If so, we should also check for its consequences, or  
 2664 at least some of them. The consequences of the original principle (1) are  
 2665 well known. We may follow here Hild (1998a), but need not do it very far.

15 In response to such examples, Briggs (2009, 71) proposes a principle of distorted reflection:  $P_0(A \mid \{\pi_1 \mid \pi_1(A) = x\}) = x - D$ , where  $D$  is a factor measuring the "expected departure from conditionalization on veridical evidence" (regarding  $A$ ). This may be a correct qualification. But again, it takes two steps at once, and we should first get clear about the single steps.

2666 Formally, the consequences of the generalized (3) are obviously analogous.  
 2667 So, the strategy in this section will be to develop the formal analogy and to  
 2668 check whether the results are also intuitively acceptable.

2669 The first thing to do, perhaps, is to unpack again what we have packed  
 2670 into the condensed abstract statement of (3). It contains, in fact, five different  
 2671 principles, depending on the temporal and the superiority relations between  
 2672  $P$  and  $P'$  (where, once and for all, each principle quantifies over all *worldly*  
 2673 propositions  $A$ ).

2674 One case is where  $P = P_0$  is the prior and  $P' = P_1$  is the posterior. If  $P_1$  is  
 2675 superior to  $P_0$ , we get:

2676 (3a)  $P_0(A \mid \{\pi_1 \mid \pi_1 = P_1\}) = P_1(A)$ , given that  $P_1$  is taken to be superior to  $P_0$ .

2677 This is our original forward-looking reflection principle (1). Since the proviso  
 2678 takes care of the main objections against the principle, there is no need to  
 2679 discuss it further.

2680 However,  $P_1$  may also be inferior to  $P_0$ . Then we get something we have not  
 2681 yet explicitly stated:

2682 (3b)  $P_0(A \mid \{\pi_1 \mid \pi_1 = P_1\}) = P_0(A)$ , given that  $P_1$  is taken to be inferior to  $P_0$ .

2683 Given tomorrow's inferior opinion, I stick to my prior opinion. E.g., today, I  
 2684 believe that I have a date with my dentist on Tuesday next week. It is only  
 2685 reasonable, then, to stick to this belief, given I am confused tomorrow and  
 2686 think the date is next Wednesday. Christensen (1991, 234) imagines an agent  
 2687 having swallowed a hefty dose of a certain drug and then being asked: "What  
 2688 do you think the probability is that you'll be able to fly in one hour, given  
 2689 that you'll then take the probability that you can fly to be .99?" He answers  
 2690 in place of the agent: "The sane answer to the above question is clearly one  
 2691 that gives a very low probability to the agent's ability to fly one hour from  
 2692 now, even on the supposition that she will at that time give it a very high  
 2693 probability" (Christensen 1991, 235). This is clearly an instance of (3b). Hence,  
 2694 Christensen may be said to have anticipated the intention of the full principle  
 2695 (3).

2696 Another case is where  $P = P_1$  is the posterior and  $P' = P_0$  is the prior. Again,  
 2697 this splits up into:

2698 (3c)  $P_1(A \mid \{\pi_0 \mid \pi_0 = P_0\}) = P_0(A)$ , given that  $P_1$  is taken to be inferior to  $P_0$ ,

2699 which is our backward-looking reflection principle (2). And into:



2700 (3d)  $P_1(A \mid \{\pi_0 \mid \pi_0 = P_0\}) = P_1(A)$ , given that  $P_1$  is taken to be superior to  $P_0$ ,

2701 which we have not yet explicitly stated, either. It says that, given my prior  
2702 opinion, I stick to my posterior opinion if it has been acquired through a  
2703 favorable change. Again, this seems to go without saying.

2704 There is, finally, the case where  $P$  and  $P'$  refer to the same time, so that  
2705  $P = P' = P_i$  ( $i = 0, 1$ ), where we may, as mentioned, define the superiority  
2706 relation either way. Thereby, we get a *synchronic reflection principle*, which is  
2707 independent of the previous diachronic principles:

2708 (3e)  $P_i(A \mid \{\pi \mid \pi = P_i\}) = P_i(A)$  ( $i = 0, 1$ ).

2709 This may look odd. But it simply says that your present opinion, whether  
2710 prior or posterior, is presently an expert for you. You presently don't know  
2711 better than you actually know. Of course, this does not preclude you from  
2712 accepting other hypothetical experts as well. Since it seems to differ from the  
2713 other principles, it has also received a separate discussion. We need not go  
2714 into it now.<sup>16</sup>

2715 What is the relation between the five parts (3a–e) of the full reflection  
2716 principle (3)? As far as I see, they are independent. As already observed,  
2717 the synchronic principle (3e) must be independent of the other diachronic  
2718 principles. Moreover, the principles (3a + d) for favorable changes and the  
2719 principles (3b + c) for unfavorable changes are independent as well, simply  
2720 because they refer to disjoint conditions. Maybe the two principles about  
2721 favorable changes are related? And likewise, those about unfavorable changes?  
2722 However, I have not discovered any relation, and I think that the five parts  
2723 (3a–e) are indeed independent. In the next section, however, I shall indicate  
2724 how things may change.

2725 Let's look at some consequences of the original reflection principle (1), or  
2726 (3a), in order to check whether their formal generalization is also intuitively  
2727 plausible. The first is the *iteration principle* already mentioned:

2728 (4a) for all worldly propositions  $A$   $P_0(A) = \Sigma \pi_1(A) \cdot P_0(\pi_1)$ , where  $P_0(\pi_1)$  is  
2729 the subject's prior auto-epistemic probability for  $\pi_1$  being her posterior,

16 Van Fraassen (1984, 248) takes it to be “uncontroversial.” However, Christensen (2007, 322–323), after calling (3a) a principle of epistemic self-respect and quoting a lot of support for it, puts forward putative counter-examples. They trade, I think, on a subtle ambiguity of the inner and outer  $P_i$  in (3a). There, the outer  $P_i$  has some (second-order) information about the inner  $P_i$ , and, thus, the two may come to diverge.

2730 where the sum is taken over all her possible posteriors  $\pi_1$  taken to be  
 2731 superior to  $P_0$ , and where  $\Sigma P_0(\pi_1) = 1$ , i.e.,  $P_0$  is sure to undergo a  
 2732 favorable change.

2733 (4a) is entailed by (3a);<sup>17</sup> for a possible reversal, see below. In other words,  
 2734 your prior opinion is always a weighted mixture of all the posterior opinions  
 2735 you may favorably reach, where the weights are given by your prior opinion  
 2736 about reaching these posteriors.<sup>18</sup> This is, I think, a deep epistemological  
 2737 insight.

2738 In the same way, the backward-looking reflection principle (2), or (3c),  
 2739 entails the following *reverse iteration principle*:

2740 (4b) for all worldly propositions  $A$   $P_1(A) = \Sigma \pi_0(A) \cdot P_1(\pi_0)$ , where the sum  
 2741 is taken over all possible priors  $\pi_0$  taken to be superior to  $P_1$ , and again  
 2742  $\Sigma P_1(\pi_0) = 1$ , i.e.,  $P_1$  is sure to have undergone an unfavorable change.

2743 The proof is analogous to the one of (4a). Is (4b) plausible? Yes. If, in your poster-  
 2744 ior  $P_1$ , you have forgotten about something, you will usually not remember,  
 2745 either, what your past opinion about that thing has been. Still, you might auto-  
 2746 epistemically wonder what your past opinion has been. And this guesswork is  
 2747 coherent only if it satisfies (4b). For instance, you cannot coherently say, “Oh,  
 2748 I have forgotten my date with the dentist; I guess it’s next Wednesday. But I  
 2749 think that yesterday I was still quite sure that it was next Tuesday.” You may  
 2750 reversely take this as support for the backward-looking reflection principle  
 2751 (2).

2752 One may think that there is a difference between (4a) and (4b). (4a), but not  
 2753 (4b), is grounded, as it were, in experience. The standard instantiation of (4a)  
 2754 is simple conditionalization: You learn exactly one member of a partition of ev-  
 2755 idential (worldly) propositions about which you have some prior expectations.  
 2756 And since your possible posteriors are just the conditionalization of your  
 2757 prior with respect to these evidential propositions, you have the very same  
 2758 expectations about these posteriors. Such a grounding is, however, entirely  
 2759 missing in the case of (4b). You may remember your past probabilities, but if

17 If we suppress the additional condition about superiority and stick to the shorter notation used in (4a), (3a) says  $P_0(A | \pi_1) = \pi_1(A)$ . The formula of the total probability says that  $P_0(A) = \Sigma P_0(A | \pi_1) \cdot P_0(\pi_1)$ . By inserting the first equation into the second, we get (4a).

18 Here, de Finetti’s heritage is particularly salient. Recall his famous representation theorem, saying that your (prior) probabilities are symmetric or exchangeable (as they should be) if and only if they are a unique mixture of all the statistical hypotheses they might converge to.

2760 you have forgotten them, your present opinion about them is mere guesswork  
 2761 without such grounding.

2762 However, the case is not as asymmetric as it seems. (4a) is not only made for  
 2763 simple conditionalization. It holds as well, e.g., for Jeffrey conditionalization,  
 2764 where learning results in some new posterior probabilities for the partition of  
 2765 evidential propositions. And then the posterior is not grounded in a specific  
 2766 evidential proposition but in your possibly vague seemings concerning this  
 2767 evidential partition. Then, however, your expectations about these seemings  
 2768 are not much better off than in the past-oriented case. So, (4a + b) is auto-  
 2769 epistemic business justified by full reflection (3). Such specific grounding is  
 2770 welcome but not required.

## 2775 **Reflection Applied to Auto-Epistemic Propositions**

2772 To check out further consequences, we must attend to the way Hild (1998a)  
 2773 presents the principles. He tacitly assumes an innocent-looking generaliza-  
 2774 tion; i.e., from the outset, he applies van Fraassen’s principle (1) also to auto-  
 2775 epistemic (and mixed) propositions  $A$ . Let’s briefly consider in this section  
 2776 what happens when we thus generalize our principles (3a–e) and drop their  
 2777 restriction to worldly propositions. Note that this also requires us to drop  
 2778 the restriction of the possible probability measures  $\pi$  to worldly propositions.  
 2779 Hence, equations like  $\pi = Q$  can now be literally and not only sloppily true.

2780 A first consequence would be that the synchronic reflection principle (3e)  
 2781 turns out to be equivalent to what Hild calls *auto-epistemic transparency*.<sup>19</sup>

$$2782 \quad (5) \quad P_i(\{\pi \mid \pi = P_i\}) = 1 \quad (i = 0, 1).^{20}$$

2783 In other words, in each second-order epistemic state, also reflecting on your  
 2784 first-order state, you know, or are sure, what your present first-order state is.  
 2785 In doxastic logic, this is sometimes called “positive introspection” or the BB  
 2786 thesis, “If you believe that  $p$ , then you believe that you believe that  $p$ ,” first  
 2787 discussed in Hintikka (1962, 123ff.), amply attacked and amply defended. Let  
 2788 us not engage in this discussion now.<sup>21</sup>

19 Christensen (2007) and Weisberg (2007) split this up into two principles called confidence and accuracy by Christensen and luminosity and transparency by Weisberg.

20 *Proof*: For the one direction, take  $A$  in (3e) to be the auto-epistemic proposition  $\{\pi \mid \pi = P_i\}$ . Reversely,  $P_i$  is identical to  $P_i$  conditional on a proposition with probability 1.

21 Besides the arguments referred to in footnote 15, Christensen (2007) casts doubt on auto-epistemic transparency by questioning that we have certain knowledge about our precise subjective probabil-

2789 Another consequence of the reflection principles (3a + e) thus extended is  
 2790 *perfect memory*:

2791 (6)  $P_1(\{\pi_0 \mid \pi_0 = P_0\}) = 1$ , given that  $P_1$  is taken to be superior to  $P_0$ .<sup>22</sup>

2792 This is a suspiciously strong consequence. However, given the extension of  
 2793 the principles to auto-epistemic propositions, one could say that whenever I  
 2794 have become uncertain about my former epistemic state, I have forgotten my  
 2795 former attitude towards some proposition. So, my uncertainty must be the  
 2796 result of an unfavorable change.

2797 By a similar proof, the extended (3c) + (3e) implies an analogous principle  
 2798 of perfect foresight. My prior knows for sure what my inferior posterior will be.  
 2799 This is obviously absurd. Hence, the auto-epistemic extension of the backward-  
 2800 looking principle (3c) must be rejected. The proof of (6) displays where the  
 2801 analogy breaks down. Inserting  $\{\pi_1 \mid \pi_1 = P_1\}$  for  $A$  in (3c) would mean that I  
 2802 would trust my former superior  $P_0$  concerning my present inferior state. But  
 2803 regarding my own present state, I am always in an optimal epistemic position,  
 2804 as confirmed by auto-epistemic transparency (5); in this respect, I need no  
 2805 lessons from my better self. My position may be inferior only with respect  
 2806 to worldly propositions (and auto-epistemic propositions referring to other  
 2807 times). This is why we must not extend (3c) to auto-epistemic propositions  
 2808 (or at least not to the simultaneous ones).

2809 The envisaged extension also helps a bit regarding the relations among the  
 2810 five parts of (3). That is, we may see now that the extended (3a + e) entails  
 2811 not only (6) but also (3d) simply because (6) says that the condition of (3d)  
 2812 has probability 1. For the same reason as before, we must not exploit this  
 2813 observation for a corresponding derivation of (3b) from (3c + e); (3b) seems  
 2814 to remain independent.

---

ities. However, this rather questions the representation of epistemic states by precise probabilities, which we have assumed at the outset of this paper. That is, I tend to assume that beliefs and epistemic states, in general, are conscious mental states (in the sense of intentional or higher-order consciousness), similar to phenomenally conscious pains. So, if I do not know my precise probabilities, I don't have them, just as I don't have pains when I am not aware of them. It is this assumption, I think, that is the motivation behind the BB thesis and its kin.

22 *Proof:* Take  $A$  in (3a) to be the auto-epistemic proposition  $\{\pi_0 \mid \pi_0 = P_0\}$ . Thus,  $P_0(\{\pi_0 \mid \pi_0 = P_0\} \mid \{\pi_1 \mid \pi_1 = P_1\}) = P_1(\{\pi_0 \mid \pi_0 = P_0\})$ , given that  $P_1$  is taken to be superior to  $P_0$ . Auto-epistemic transparency (5) says that  $P_0(\{\pi_0 \mid \pi_0 = P_0\}) = 1$ . Hence,  $P_0(\{\pi_0 \mid \pi_0 = P_0\} \mid \{\pi_1 \mid \pi_1 = P_1\}) = 1$  as well. So, finally,  $P_1(\{\pi_0 \mid \pi_0 = P_0\}) = 1$ , given that  $P_1$  is taken to be superior to  $P_0$  (cf. Hild 1998a, 353).

2815 Moreover, we may note that Hild's extension strengthens the relation be-  
 2816 tween reflection and iteration. We observed already that (3a) and (3c), respec-  
 2817 tively, entail (4a) and (4b). With the extension, we may reverse the entailment:  
 2818 given auto-epistemic transparency (5) or the equivalent synchronic reflection  
 2819 principle (3e), the iteration principle (4a) implies the forward-looking reflec-  
 2820 tion principle (3a).<sup>23</sup> Thus, under the same assumptions, the reverse iteration  
 2821 principle (4b) entails the backward-looking (3c), as seems unobjectionable.

2822 Finally, this extension helps us to provide a formal proof of my informally  
 2823 justified claim in section 2 that the forward- and backward-looking principles  
 2824 (1) and (2) can apply simultaneously only in the case of non-change. Given  
 2825 that (1) and (2) apply also to auto-epistemic propositions, we have:

2826 (7) Given auto-transparency (5), if (a)  $P_0(A \mid \{\pi_1 \mid \pi_1 = P_1\}) = P_1(A)$  and (b)  
 2827  $P_1(A \mid \{\pi_0 \mid \pi_0 = P_0\}) = P_0(A)$  hold for all propositions  $A$ , then  $P_0 = P_1$ .<sup>24</sup>

2828 In conclusion, we have found a slight divergence among our principles in this  
 2829 extension, a divergence we could justify. In the main, however, the parallel  
 2830 between the forward-looking and the backward-looking perspective and thus  
 2831 between the parts of the full reflection principle (3) stands. We have not  
 2832 discovered any incoherence.

## 2833 6 The Full Reflection Principle and the Expert Principle

2834 In section 1, I pointed already to the expert principle, which is the most  
 2835 common generalization of van Fraassen's principle (1). It may seem that the  
 2836 full reflection principle (3) is just another special case of the expert principle.<sup>25</sup>  
 2837 Yes, almost. At least, this holds for the backward-looking principle (2). Here,  
 2838 my better-informed past self may be taken to be an expert for my present  
 2839 forgetful self. However, not all cases of (3) are special cases of the expert  
 2840 principle. Let  $P$  in (3) be my probability measure and  $P'$  that of my neighbor.  
 2841 When I take my neighbor to be better informed, to be in a superior epistemic

23 *Proof:* We have  $P_0(A \text{ and } \{\pi_1 \mid \pi_1 = P_1\}) = \sum \pi(A \text{ and } \{\pi_1 \mid \pi_1 = P_1\}) \cdot P_0(\pi)$  (by the extended (4a), where the sum is taken over all possible posteriors  $\pi = P_1(A) \cdot P_0(\{\pi_1 \mid \pi_1 = P_1\})$  (by auto-transparency, because  $\pi(\{\pi_1 \mid \pi_1 = P_1\}) = 1$  only for  $\pi = P_1$ , and otherwise = 0) (cf. Hild 1998a, 354).

24 *Proof:* We have  $P_1(A) = P_0(A \mid \pi_1 = P_1)$  (due to (a)) =  $P_0(A \mid \pi_1 = P_1, \pi_0 = P_0)$  (due to auto-transparency (5)) =  $P_1(A \mid \pi_0 = P_0)$  (by applying (a) conditional on  $\pi_0 = P_0$ ) =  $P_0(A)$  (due to (b)).

25 I am grateful to a reviewer for raising this issue.

2842 position (concerning a certain field), then I listen to her (in the sense of  
 2843 obeying (3)); these are the cases (3a) and (3c), where the temporal relation  
 2844 between me and my neighbor is irrelevant. But when I take her to be less well  
 2845 informed or in an inferior epistemic position, I do not listen to her; these are  
 2846 the cases (3b) and (3d). This seems to go without saying. Christensen (2000,  
 2847 358) takes this for granted, too. Strictly speaking, though, it is not part of the  
 2848 expert principle, which says only how to deal with people taken to be at least  
 2849 as well-informed.

2850 Of course, it would be no problem to pair the expert principle with a “non-  
 2851 expert principle,” saying that, rationally, we don’t listen to persons we take  
 2852 to be in an inferior epistemic position. However, this would still leave us  
 2853 with very incomplete principles. The no-neutrality condition—which was  
 2854 plausible in the intrasubjective case, at least when we can divide up complex  
 2855 epistemic changes into unidirectional steps—has no analog with respect to  
 2856 experts. Most of my fellow humans are neither better nor less well-informed  
 2857 than me; their epistemic state is just incomparable to mine. And then both  
 2858 the expert and the non-expert principle are silent. This is not an objection. It  
 2859 is hard to give any recommendations for the incomparable cases. But it is our  
 2860 daily business to somehow deal with them.

2861 One may think<sup>26</sup> that we can apply the treatment of intrasubjective incom-  
 2862 parabilities suggested in section 3 also to the interpersonal case of experts.  
 2863 However, this is not so easy. In the intrapersonal case, we had to refer, it  
 2864 seemed, to counterfactual epistemic states that the subject would be in had  
 2865 certain unfavorable changes not occurred. This might be manageable. In the  
 2866 case of my incomparable neighbor, however, the corresponding counterfactual  
 2867 question would be which favorable changes he would have to undergo  
 2868 and which unfavorable changes to avoid till I could acknowledge him to be  
 2869 an expert, i.e., to be in an equal or superior epistemic position concerning the  
 2870 issue at hand. This is a much more sweeping and indeterminate counterfactual  
 2871 question. The point is that the superiority and inferiority of epistemic  
 2872 positions are clearly assessable on the basis of intrasubjective favorable or  
 2873 unfavorable changes. But it is very hard to assess as such, as an interpersonal  
 2874 comparison would require.

2875 So, we have a tension here. While van Fraassen’s principle (1) was clearly a  
 2876 special case of the expert principle, the subsumption of the full principle (3)  
 2877 is at least doubtful. In fact, such a subsumption was not intended in the begin-

---

26 Suggested by the same reviewer.

2878 ning. Gaifman (1986) proposed the expert principle as a formal generalization  
 2879 of the reflection principle, not as a substitute for the latter’s epistemological  
 2880 role. As such it operates only in so-called time-slice epistemology (Moss 2015)  
 2881 or time-slice rationality (Hedden 2015). A basic assumption of this approach  
 2882 is called impartiality:<sup>27</sup> “In determining how you rationally ought to be at a  
 2883 time, your beliefs about what attitudes you have at other times play the same  
 2884 role as your beliefs about what attitudes other people have” (Hedden 2015, 9).  
 2885 If so, it is clear that reflection principles are unnecessarily restrictive and that  
 2886 the expert principle completely takes over the epistemological role of the full  
 2887 reflection principle (given an additional “non-expert principle”).

2888 In his defense of time-slice rationality, Hedden (2015, chaps. 8–9) makes  
 2889 crucial use of an assumption called uniqueness (by Feldman 2007): “Given a  
 2890 body of total evidence, there is a unique doxastic state that it is rational to be  
 2891 in” (Hedden 2015, 130). This reveals an entirely different picture of normative  
 2892 epistemology than the one pursued here. It is that there is a unique prior—the  
 2893 *ur*-prior, as it were—and then all epistemic change is due to a change of the  
 2894 body of total evidence.<sup>28</sup> Given this, we do not need any diachronic rules for  
 2895 epistemic change. We always refer back to the *ur*-prior and then consider  
 2896 how the given body of total evidence operates on it. All change is in that  
 2897 body and only there; the body may get larger (through learning) or smaller  
 2898 (through forgetting).<sup>29</sup> My fellow humans are in the very same situation. They  
 2899 rationally proceed from the very same *ur*-prior, and they differ from me only  
 2900 in their total evidence. This is why they count just as much as my future  
 2901 or past epistemic states. Or rather, nobody counts; only the bodies of total  
 2902 evidence count. Certainly, this would simplify our epistemological business  
 2903 considerably.

2904 The assumption of uniqueness also makes the notion of an expert very easy.  
 2905 Among rational subjects, *a* is an expert for *b* simply if *a* has at least as much  
 2906 evidence as *b*. If *a*’s and *b*’s evidence only overlap, they are in incomparable  
 2907 states. But joining their evidence would result in an expert for both. With

27 The other basic assumption is synchronicity: “All rationality requirements are synchronic” (Hedden 2015, 9). See also Moss (2015, 177) for a statement of the two principles. The principle of impartiality seems to have been stated first by Christensen (2000, 363–364). Moss (2015, 178) and Hedden (2015, 56) happily observe that van Fraassen’s reflection principle satisfies the assumption of synchronicity insofar as it speaks only about  $P_0$  and its conjectures about  $P_1$ . Of course, this observation carries over to the other versions.

28 Or the total evidence need not appeal to any prior at all. But then the *ur*-prior is unique as well, namely, empty.

29 This picture also motivates Titelbaum’s (2013) framework of gain and loss of certainties.

2908 this easy notion of an expert, the above idea of paralleling intersubjective  
 2909 incomparabilities with intrasubjective incomparable changes might be less  
 2910 problematic.

2911 The alternative is to deny uniqueness. Hedden (2015, 129) calls this per-  
 2912 missiveness. But what is the dialectic situation here? Pace Hedden, it is not  
 2913 that one must defend uniqueness and the other permissiveness (by showing  
 2914 two *ur*-priors to be equally rationally acceptable). In my view, the burden of  
 2915 proof lies with the defender of uniqueness. And a proof should constructively  
 2916 indicate what the unique *ur*-prior looks like. The literature is not so promising.  
 2917 The only positive attempt I know of is objective Bayesianism, as proposed  
 2918 by Williamson (2005).<sup>30</sup> By contrast, I tend to take centuries of skepticism to  
 2919 suggest that such a proof will fail.<sup>31</sup> Obviously, this is too big an issue to be  
 2920 discussed now.<sup>32</sup> The point is only this: In the absence of such a proof, we  
 2921 should not proceed from the assumption of uniqueness. A positive defense of  
 2922 permissiveness is not really required.<sup>33</sup>

2923 Kelly (2014) usefully distinguishes statements of uniqueness that have  
 2924 interpersonal import, as he calls it, from those that have not. Intrasubjective  
 2925 uniqueness only requires that, given my background or my personal *ur*-prior,  
 2926 I have only one rational way of responding to the evidence. I have no quarrel  
 2927 with this. However, uniqueness with interpersonal import requires that there  
 2928 is only one and the same rational way for everybody to respond to the evidence.  
 2929 This is the version intended and critically discussed above.

2930 The question, then, is how to pursue normative epistemology without the  
 2931 assumption of uniqueness. Just in the way it is done traditionally, and here

---

30 To be precise, Williamson does not need to refer to an *ur*-prior. He rather proposes a unique way of responding to any given total body of evidence, which does not depend on an underlying *ur*-prior. Rather, the *ur*-prior would be the response to the empty evidence.

31 Hedden (2015, 134) himself (as well as Kelly 2014, 309) points to the alleged failure of Carnap's (1962) project of inductive logic, which also started out searching for the unique prior. I should add, though, that in chap. 8, Hedden admits that uniqueness may force one to allow for indeterminate and/or imprecise probabilities. However, considering other epistemic formats shifts the discussion still further. I have explained why I focus here on precise probabilities only.

32 I admit that the issue can also be discussed in the abstract without constructive proposals for an *ur*-prior. See, e.g., the exchange between White (2014) and Kelly (2014). It is clear that my sympathies lie here with Kelly. However, he is still too obliging, I find; he does not raise the point about the burden of proof.

33 Recall also this: In the final section of his (1980), Lewis discovers a tension between his Principal Principle and Humean Supervenience. He considers resolving the tension by assuming what is now called uniqueness. But he shies away from this solution, which he finds "not very easy to believe." As is well known, he modified the Principal Principle later on.



2932 as well, namely, by stating synchronic principles of epistemic rationality and  
 2933 diachronic principles. The latter can only refer to a subject's prior and posterior  
 2934 and a piece of total evidence in between, but *not* to an *ur*-prior and a body of  
 2935 total evidence reaching back to the indefinite time of the *ur*-prior.<sup>34</sup> In this  
 2936 conception, the reflection principles, as discussed here, have a natural place,  
 2937 and the goal of stating a complete dynamic is important, while atemporal  
 2938 expert principles do not directly add to it and need not be complete. Given  
 2939 uniqueness, we can also distinguish inferior and superior epistemic positions  
 2940 simply by looking at the size of the bodies of total evidence underlying them.  
 2941 However, when looking at epistemic dynamics in the traditional way, then, as  
 2942 argued, we must also classify single changes as favorable and unfavorable (and  
 2943 if we cannot do this objectively, we leave it to the subject herself, as proposed  
 2944 here). As mentioned, favorable and unfavorable changes do not only consist  
 2945 in gaining and losing evidence or certainties; they may take various other  
 2946 forms not easily subsumed under the picture motivated by uniqueness. In the  
 2947 perspective pursued here, expert principles become relevant only because we  
 2948 take listening to experts to induce favorable change, unlike listening to non-  
 2949 experts. This is why I think that the reflection principles have an independent  
 2950 value. They can be substantially subsumed under the expert principles only  
 2951 within a questionable epistemological picture.\*

2952 Wolfgang Spohn

2953  0000-0002-3213-8907

2954 Department of Philosophy, University of Konstanz  
 2955 wolfgang.spohn@uni-konstanz.de

## 2956 References

- 2957 ARNTZENIUS, Frank. 2003. "Some Problems for Conditionalization and Reflection."  
 2958 *The Journal of Philosophy* 100(7): 356–370, doi:10.5840/jphil2003100729.

---

34 When we model a dynamic process, in physics, meteorology, or wherever, we do it in the form of a law saying how one state of the system modeled changes into the subsequent state, possibly under the influence of external factors, and we can do this in discrete or in continuous time. So, this is also the natural format for normative epistemology as well, where we try to say what a rational epistemic dynamic should look like.

\* I am deeply indebted to three reviewers, one from the *Journal of Philosophy* and two from *Dialectica*. They considerably helped me in various respects. I also acknowledge support by the German Research Foundation, Grant EXC 2064/1, No. 390727645, and Grant SP279/21-1, No. 420094936.

- 2959 BINKLEY, Robert W. 1968. "The Surprise Examination in Modal Logic." *The Journal of*  
 2960 *Philosophy* 65(5): 127–136, doi:10.2307/2024556.
- 2961 BRIGGS, Rachael A. 2009. "Distorted Reflection." *The Philosophical Review* 118(1):  
 2962 59–85, doi:10.1215/00318108-2008-029.
- 2963 CARNAP, Rudolf. 1950. *Logical Foundations of Probability*. Chicago, Illinois: University  
 2964 of Chicago Press. Second edition: Carnap (1962).
- 2965 —. 1962. *Logical Foundations of Probability*. 2nd ed. Chicago, Illinois: University of  
 2966 Chicago Press. First edition: Carnap (1950).
- 2967 CHRISTENSEN, David. 1991. "Clever Bookies and Coherent Beliefs." *The Philosophical*  
 2968 *Review* 100(2): 229–247, doi:10.2307/2185301.
- 2969 —. 2000. "Diachronic Coherence Versus Epistemic Impartiality." *The Philosophical*  
 2970 *Review* 109(3): 349–371, doi:10.1215/00318108-109-3-349.
- 2971 —. 2007. "Epistemic Self-Respect." *Proceedings of the Aristotelian Society* 107: 319–337,  
 2972 doi:10.1111/j.1467-9264.2007.00224.x.
- 2973 —. 2010a. "Higher-Order Evidence." *Philosophy and Phenomenological Research* 81(1):  
 2974 185–215, doi:10.1111/j.1933-1592.2010.00366.x.
- 2975 —. 2010b. "Rational Reflection." in *Philosophical Perspectives 24: Epistemology*, edited  
 2976 by John HAWTHORNE, pp. 121–140. Hoboken, New Jersey: John Wiley and Sons,  
 2977 Inc., doi:10.1111/j.1520-8583.2010.00187.x.
- 2978 FELDMAN, Fred. 2004. *Pleasure and the Good Life: Concerning the Nature, Varieties,*  
 2979 *and Plausibility of Hedonism*. Oxford: Oxford University Press, doi:10.1093/0199  
 2980 26516X.001.0001.
- 2981 —. 2007. "Précis of Feldman (2004)." *Philosophical Studies* 136(3): 405–408, doi:10.100  
 2982 7/s11098-006-9039-y.
- 2983 DE FINETTI, Bruno. 1937. "La prévision, ses lois logiques, ses sources subjectives." *Annales de l'Institut Henri Poincaré* 7: 1–68. Translated as de Finetti (1964).
- 2984 —. 1964. "Foresight: Its Logical Laws, Its Subjective Sources." in *Studies in Subjective*  
 2985 *Probability*, edited by Henry E. KYBURG Jr. and Howard E. SMOKLER, pp. 93–158.  
 2986 Hoboken, New Jersey: John Wiley and Sons, Inc. Translation of de Finetti (1937),  
 2987 reprinted in Kyburg and Smokler (1980, 53–118).
- 2988 VAN FRAASSEN, Bas C. 1984. "Belief and the Will." *The Journal of Philosophy* 81(5):  
 2989 235–256, doi:10.2307/2026388.
- 2990 —. 1995. "Belief and the Problem of Ulysses and the Sirens." *Philosophical Studies*  
 2991 77(1): 7–37, doi:10.1007/bf00996309.
- 2992 GAIFMAN, Haim. 1986. "A Theory of Higher Order Probabilities." in TARK 1986. *The-*  
 2993 *oretical Aspects of Reasoning about Knowledge: Proceedings of the First Confer-*  
 2994 *ence*, edited by Joseph Y. HALPERN, pp. 275–292. San Francisco, California: Mor-  
 2995 gan Kaufmann Publishers. Reprinted in Skyrms and Harper (1988, 191–219),  
 2996 doi:10.1016/B978-0-934613-04-0.50023-5.
- 2997 GIGERENZER, Gerd. 2007. *Gut Feelings: Short Cuts to Better Decision Making*. London:  
 2998 Penguin Books.
- 2999

- 3000 GOLDSTEIN, Michael. 1983. "The Prevision of a Prevision." *Journal of the American*  
3001 *Statistical Association* 78(384): 817–819, doi:[10.2307/2288190](https://doi.org/10.2307/2288190).
- 3002 HEDDEN, Brian. 2015. *Reasons Without Persons: Rationality, Identity, and Time*. Oxford:  
3003 Oxford University Press, doi:[10.1093/acprof:oso/9780198732594.001.0001](https://doi.org/10.1093/acprof:oso/9780198732594.001.0001).
- 3004 HILD, Matthias. 1998a. "Auto-Epistemology and Updating." *Philosophical Studies* 92(3):  
3005 321–361, doi:[10.1023/a:1004229808144](https://doi.org/10.1023/a:1004229808144).
- 3006 —. 1998b. "The Coherence Argument Against Conditionalization." *Synthese* 115(2):  
3007 229–258, doi:[10.1023/a:1005082908147](https://doi.org/10.1023/a:1005082908147).
- 3008 HINTIKKA, Jaakko. 1962. *Knowledge and Belief: An Introduction to the Logic of the Two*  
3009 *Notions*. Ithaca, New York: Cornell University Press.
- 3010 KELLY, Thomas. 2014. "Evidence Can Be Permissive." in *Contemporary Debates in*  
3011 *Epistemology*, edited by Matthias STEUP, John TURRI, and Ernest SOSA, 2nd ed., pp.  
3012 298–311. Contemporary Debates in Philosophy n. 3. Malden, Massachusetts: Wiley-  
3013 Blackwell. First edition: Sosa and Steup (2005), doi:[10.1002/9781394260744.ch12](https://doi.org/10.1002/9781394260744.ch12).
- 3014 KYBURG, Henry E., Jr. and SMOKLER, Howard E., eds. 1980. *Studies in Subjective*  
3015 *Probability*. 2nd ed. Huntington, New York: Robert E. Krieger Publishing Co.
- 3016 LEWIS, David. 1980. "A Subjectivist's Guide to Objective Chance." in *Studies in Induc-*  
3017 *tive Logic and Probability: Volume II*, edited by Richard C. JEFFREY, pp. 263–294.  
3018 Berkeley, California: University of California Press. Reprinted, with a postscript  
3019 (Lewis 1986b), in Lewis (1986a, 83–113), doi:[10.1525/9780520318328-009](https://doi.org/10.1525/9780520318328-009).
- 3020 —. 1986a. *Philosophical Papers, Volume 2*. Oxford: Oxford University Press, doi:[10.1093/0195036468.001.0001](https://doi.org/10.1093/0195036468.001.0001).
- 3021 —. 1986b. "Postscript to Lewis (1980)." in *Philosophical Papers, Volume 2*, pp. 114–132.  
3022 Oxford: Oxford University Press, doi:[10.1093/0195036468.003.0004](https://doi.org/10.1093/0195036468.003.0004).
- 3023 MAHTANI, Anna. 2017. "Deference, Respect and Intensionality." *Philosophical Studies*  
3024 174(1): 163–183, doi:[10.1007/s11098-016-0675-6](https://doi.org/10.1007/s11098-016-0675-6).
- 3025 MOSS, Sarah. 2015. "Time-Slice Epistemology and Action Under Indeterminacy." in  
3026 *Oxford Studies in Epistemology*, volume V, edited by Tamar Szabó GENDLER and  
3027 John HAWTHORNE, pp. 172–194. Oxford: Oxford University Press, doi:[10.1093/acprof:oso/9780198722762.003.0006](https://doi.org/10.1093/acprof:oso/9780198722762.003.0006).
- 3028 ROTT, Hans. 2009. "Shifting Priorities: Simple Representations for Twenty-Seven Iter-
- 3029 ated Theory Change Operators." in *Towards Mathematical Philosophy: Papers from*  
3030 *the Studia Logica Conference Trends in Logic IV*, edited by David C. MAKINSON,  
3031 Jacek MALINOWSKI, and Heinrich Theodor WANSING, pp. 269–296. Trends in  
3032 Logic n. 28. Dordrecht: Springer, doi:[10.1007/978-1-4020-9084-4\\_14](https://doi.org/10.1007/978-1-4020-9084-4_14).
- 3033 SCHOENFIELD, Miriam. 2012. "Chilling Out on Epistemic Rationality: A Defense of  
3034 Imprecise Credences (and Other Imprecise Doxastic Attitudes)." *Philosophical*  
3035 *Studies* 158(2): 197–219, doi:[10.1007/s11098-012-9886-7](https://doi.org/10.1007/s11098-012-9886-7).
- 3036 SKYRMS, Brian. 1980. *Causal Necessity: A Pragmatic Investigation of the Necessity of*  
3037 *Laws*. New Haven, Connecticut: Yale University Press.
- 3038
- 3039

- 3040 SKYRMS, Brian and HARPER, William L., eds. 1988. *Causation, Chance, and Credence:*  
 3041 *Proceedings of the Irvine Conference on Probability and Causation, Volume 1.* The  
 3042 University of Western Ontario Series in Philosophy of Science n. 41. Dordrecht:  
 3043 Kluwer Academic Publishers, doi:10.1007/978-94-009-2863-3.
- 3044 SOSA, Ernest and STEUP, Matthias, eds. 2005. *Contemporary Debates in Epistemology.*  
 3045 1st ed. Contemporary Debates in Philosophy n. 3. Boston, Massachusetts: Blackwell  
 3046 Publishers. Second edition: Steup, Turri and Sosa (2014).
- 3047 SPOHN, Wolfgang. 1978. *Grundlagen der Entscheidungstheorie.* Kronberg/Ts.: Scriptor  
 3048 Verlag, urn:nbn:de:bsz:352-135468.
- 3049 —. 1988. “Ordinal Conditional Functions: A Dynamic Theory of Epistemic States.”  
 3050 in *Causation in Decision, Belief Change, and Statistics. Proceedings of the Irvine*  
 3051 *Conference on Probability and Causation, vol. 2*, edited by William L. HARPER  
 3052 and Brian SKYRMS, pp. 105–134. The University of Western Ontario Series in  
 3053 Philosophy of Science n. 42. Dordrecht: Kluwer Academic Publishers, doi:10.100  
 3054 7/978-94-009-2865-7\_6.
- 3055 —. 2010. “Chance and Necessity: From Humean Supervenience to Humean Projection.”  
 3056 in *The Place of Probability in Science: In Honor of Ellery Eells (1953–2006)*, edited  
 3057 by James H. FETZER, pp. 101–131. Boston Studies in the Philosophy of Science  
 3058 n. 284. Dordrecht: Springer, doi:10.1007/978-90-481-3615-5\_5.
- 3059 —. 2017. “The Epistemology and Auto-Epistemology of Temporal Self-Location and  
 3060 Forgetfulness.” *Ergo* 4(13): 359–418, doi:10.3998/ergo.12405314.0004.013.
- 3061 STEUP, Matthias, TURRI, John and SOSA, Ernest, eds. 2014. *Contemporary Debates in*  
 3062 *Epistemology.* 2nd ed. Contemporary Debates in Philosophy n. 3. Malden, Mas-  
 3063 sachusetts: Wiley-Blackwell. First edition: Sosa and Steup (2005), doi:10.1002/97  
 3064 81394260744.
- 3065 TALBOTT, William J. 1991. “Two Principles of Bayesian Epistemology.” *Philosophical*  
 3066 *Studies* 62(2): 135–150, doi:10.1007/bf00419049.
- 3067 TITELBAUM, Michael G. 2013. *Quitting Certainties: A Bayesian Framework Modeling*  
 3068 *Degrees of Belief.* Oxford: Oxford University Press, doi:10.1093/acprof:oso/978019  
 3069 9658305.001.0001.
- 3070 WEISBERG, Jonathan. 2007. “Conditionalization, Reflection, and Self-Knowledge.”  
 3071 *Philosophical Studies* 135(2): 179–197, doi:10.1007/s11098-007-9073-4.
- 3072 WHITE, Roger. 2014. “Evidence Cannot Be Permissive.” in *Contemporary Debates in*  
 3073 *Epistemology*, edited by Matthias STEUP, John TURRI, and Ernest SOSA, 2nd ed., pp.  
 3074 312–323. Contemporary Debates in Philosophy n. 3. Malden, Massachusetts: Wiley-  
 3075 Blackwell. First edition: Sosa and Steup (2005), doi:10.1002/9781394260744.ch12.
- 3076 WILLIAMSON, Jon. 2005. *Bayesian Nets and Causality: Philosophical and Computa-*  
 3077 *tional Foundations.* Oxford: Oxford University Press, doi:10.1093/acprof:oso/9780  
 3078 198530794.001.0001.

# Our Naïve Representation of Time and of the Open Future

BATOUL HODROJ, ANDREW JAMES LATHAM & KRISTIE  
MILLER

It's generally thought that we naïvely or pre-theoretically represent the future to be open. While philosophers have modelled future openness in different ways, it's unclear which, if any, captures our naïve sense that the future is open. In this paper, we focus on just one way the future might count as being open: by being nomically open. We empirically investigate whether our naïve representation of the future as open is partly constituted by representing the future as nomically open. We also investigate the connection between our naïve representation of the future as open and our representation of time. One of the purported advantages of the growing block theory of time is that it captures our naïve sense that the future is open and the past closed. We investigate whether there is an explanatory connection between people representing the future as being nomically open and representing our world as a growing block. We also reflect on the implications of our findings for theorising about future openness and temporal ontology.

It's often thought that our intuitive or pre-reflective view of the world is one in which, in some sense or other, the future is open.<sup>1</sup> It has also been thought that our intuitive, pre-reflective, or folk view of the world is one in which the totality of our world grows as new being comes into existence in the present moment and then becomes past as yet more being comes into existence.<sup>2</sup> This latter view is the view that our world is a *growing block*.<sup>3</sup>

- 
- 1 Callender (2017) takes this to be part of the manifest image; Ismael (2012) likewise.
  - 2 See Forbes (2016). Latham, Miller and Norton (2021b) confirmed empirically that, of the ~70% of people who are temporal dynamists, the most popular view is the growing block view.
  - 3 Defenders of this view include Broad (1923, 1938), Forbes (2016), Correia and Rosenkranz (2018), Tooley (1997), and Forrest (2004).

3102 In what follows, rather than talking about pre-reflective or folk views, we  
3103 will talk of *naïve representations* of the world. As we will understand them,  
3104 naïve representations are contentful mental states, i.e., representations of  
3105 various aspects of our world that are not informed by (or, at least, are largely  
3106 not the product of engagement with) current science or philosophy. These  
3107 are folk views, folk theories, or folk models of aspects of the world. These  
3108 representations may be tacit in the sense that the people whose representa-  
3109 tions they are may not be able to specify the content of the representation  
3110 when asked. Nevertheless, we take it that these representations guide people's  
3111 behaviours (linguistic and otherwise) and that we can probe their content by  
3112 giving people tasks that require them to use those representations.

3113 We are interested in two sorts of naïve representations. The first is our  
3114 *naïve representation of the future*; the second is our *naïve representation of*  
3115 *time*. Ultimately, we will be interested in whether these representations are  
3116 connected.

3117 We will take the claim that our pre-reflective view of the world is one in  
3118 which the future is open to be the claim that we naïvely represent the future  
3119 as open. Philosophers have offered various accounts of the open future. In  
3120 fact, we can (and should) distinguish at least two rather different projects with  
3121 which philosophers are engaged. The first of these aims to model the open  
3122 future. On one natural interpretation of such a project, which we will call *the*  
3123 *capturing project*, the aim is to work out which model of, or theory of, the open  
3124 future is the one that best captures our intuitive sense that the future is open.  
3125 As we construe this project, the aim is to offer a model of the open future that  
3126 best captures our naïve representation of future openness. The second project,  
3127 which we will call *the explanatory project*, focuses on explaining various “open  
3128 future” practices (conceived of very broadly) and attempts to explain why it is  
3129 that we have such practices and what it is about our world that grounds our  
3130 having such practices. These practices might include (but not be limited to)  
3131 practices of deliberating about the future but not the past; taking ourselves to  
3132 be able to causally intervene on the future but not the past; having a certain  
3133 kind of phenomenology in which the future feels, or seems, to us to be open  
3134 in the way the past does not; taking ourselves to have a kind of access to past  
3135 states that we do not have to future ones; and so on.

3136 These two projects might be connected or not. It might be that what explains  
3137 why we have the open future practices we do is the very thing that, in fact,  
3138 captures our naïve representation of the future. In that case, we will say that  
3139 our naïve representation of the future is *vindicated*. Alternatively, it could

3140 be that what explains our open future practices does not capture our naïve  
3141 representation of the future as open. To see this, consider several of the views  
3142 philosophers have put forward as models of the open future and suppose  
3143 these are claims about our naïve representation of future openness.

3144 The first view models future openness in terms of *alethic openness*. On  
3145 this view, our naïve representation of the open future consists in, or at least  
3146 includes, our representing that (some or all) future-tensed contingent state-  
3147 ments fail to take a determinate truth-value (see, for instance, [Markosian](#)  
3148 [1995](#); [Williams 2008](#) (unpublished); [MacFarlane 2003](#); and [Tooley 1997](#)). The  
3149 second of these is *epistemic openness*. On this view, our naïve representation  
3150 of the future being open consists in, or at least includes, our representing  
3151 that we have epistemic access to the future only by making predictions and  
3152 forming intentions and not by having records of what will happen (see, for  
3153 instance, [Lewis 1979](#)). The third is *nommic openness*. On this view, our naïve  
3154 representation of the future being open consists in, or at least includes, our  
3155 representing that future-directed indeterminism is true. There are multiple  
3156 ways the future could go, consistent with how it has already gone ([Belnap](#)  
3157 [1992, 2005](#); [MacFarlane 2003, 2008](#); and [McCall 1994](#)).

3158 It could be that our naïve representation of the future as open consists in our  
3159 representing the future as being open in some or all of these ways.<sup>4</sup> Suppose  
3160 it were to turn out that our naïve representation of future openness consists  
3161 entirely in representing the future to be alethically open. Suppose, however,  
3162 that our world is not, in fact, alethically open. Still, *something* explains why  
3163 we have the open future practices that we do. It might be that the fact that  
3164 there is an epistemic asymmetry between past and future is what explains  
3165 our having these practices. It might even be that the world being this way  
3166 legitimises or makes those practices rationally permissible (or obligatory).  
3167 Still, it will turn out that what explains our having the open future practices  
3168 we do does not *vindicate* our naïve representation of the future as open.

3169 This paper will have nothing to say about why we have the open future  
3170 practices we do. We set aside the explanatory project and focus entirely on  
3171 the question of what our naïve representation of future openness consists in.  
3172 This is a vital first step if we are interested in the question of whether what it  
3173 is that explains our practices (whatever that might be) vindicates our naïve  
3174 representation of the future as open.

---

4 This is not to say that these are the only such ways. For a discussion of the ways in which we could model openness, see [Torre \(2011\)](#) and [Markosian \(1995\)](#).

3175 Some work in this area has already been undertaken. Previous research by  
3176 Hodroj et al. (2023) suggests that our naïve representation of the future as  
3177 open, at least partly, consists in our representing the future to be alethically  
3178 open. So, in this paper, we focus on nomic openness. We will suppose that a  
3179 world is nomicly open just in case that world is future-wise indeterministic.  
3180 That is, a world,  $w$ , is nomicly open just in case for any time  $t$  in  $w$ , it is  
3181 not the case that a complete specification of the way the world is at  $t$ , in  
3182 conjunction with the laws of nature of  $w$ , logically entails the way the world  
3183 is at all times later than  $t$ . This leaves open that  $w$  may or may not be past-wise  
3184 nomicly open: that is, whether the way the world is at  $t$ , in conjunction with  
3185 the laws of nature, logically entails the way the world is at all times earlier  
3186 than  $t$ . Then, we are interested in whether our naïve representation of the  
3187 future involves our representing the future to be nomicly open.

3188 We are also interested in the connection between our naïve representation  
3189 of the future as open and our naïve representation of the temporal dimension.  
3190 That is because it has been suggested that part of what explains why the  
3191 growing block theory is intuitively plausible is that we naïvely represent the  
3192 future as open, and the growing block theory better captures, or better accords  
3193 with, this (see, for instance, [Briggs and Forbes 2012](#); [Forbes 2016](#); [Grandjean  
3194 2021, 2022](#); and [Correia and Rosenkranz 2018](#)).

3195 According to the growing block model of time, past events and objects exist,  
3196 but future ones do not. There is a set of events that are objectively present, and  
3197 these are the events that sit at the end of the block looking out into the non-  
3198 existent future. The temporal passage consists in the coming into existence  
3199 of a new being on the edge of reality, where that new being becomes the  
3200 objective present until more being comes to exist (at which point it becomes  
3201 part of the objective past). Hence, the growing block theory is a version of  
3202 the A-theory on which there exists robust temporal passage: there is a fact of  
3203 the matter as to which events are present and which those are changes. By  
3204 contrast, the block universe theory is a version of the B-theory. On this view,  
3205 past, present, and future events/objects exist on a four-dimensional manifold  
3206 and bear unchanging relations of earlier-than, later-than, and simultaneous-  
3207 with to one another.<sup>5</sup> None of these events is singled out as objectively present,  
3208 and so time does not robustly pass since there is no change in which events  
3209 are objectively present.

---

5 This, of course, is also true of the moving spotlight theory, which is a version of dynamism. However, on that view, unlike the block universe view, there is a single set of events singled out as objectively present.



3210 Unlike other models of time, the growing block theory has a built-in asym-  
3211 metry between the past and the future. The past exists and is located some-  
3212 where in space-time, whereas the future is yet to happen and does not exist.  
3213 By contrast, presentism holds that *neither* the future nor the past exists, and  
3214 the block universe theory holds that *both* the future and the past exist. The  
3215 moving spotlight theory also holds that both future and past exist but holds  
3216 that some events are objectively present (namely, those on which the spotlight  
3217 of presentness shines, as it were) and that which events those are changes as  
3218 the present moves.<sup>6</sup>

3219 This asymmetry has been hypothesised to better capture people's intuitive  
3220 sense that the future is open and the past is closed than do views that lack  
3221 this asymmetry.<sup>7</sup>

3222 Following Latham, Miller and Norton (2021b), we take a naïve represen-  
3223 tation of time to be a (probably tacit) representation of time and temporal  
3224 ontology in our world. People's naïve representation of time might be closer  
3225 to one or another of the models of time that philosophers engage with.

3226 Following Hodroj et al. (2023), we can distinguish three aspects of the idea  
3227 that the growing block theory better accommodates people's intuitive sense  
3228 that the future is open.

3229 First, according to the *vindication claim*, our naïve representation of future  
3230 openness has a content that is vindicated if our world is a growing block.  
3231 The narrow version of the vindication claim that will be of interest to us  
3232 in this paper is the claim that our naïve representation of future openness  
3233 has a content that is vindicated if our world is a growing block and is not  
3234 vindicated if our world is a block universe. Henceforth, we will call this the  
3235 *narrow vindication claim*.

3236 One might be particularly interested in the narrow vindication claim if one  
3237 thinks that if the growing block vindicates our naïve representation of the  
3238 open future and the block universe view does not, this gives us a reason (albeit  
3239 defeasible) to prefer the former over the latter.

3240 Second, according to the *reason claim*, people believe, perhaps tacitly, that  
3241 the fact that a world has an open future is a reason to think that that world is  
3242 a growing block world rather than a block universe world.

---

6 For empirical research into people's naïve views of time, see Latham, Miller and Norton (2021a).

7 Something that Grandjean (2021, 2022) and Correia and Rosenkranz (2018) point to.

3243 Third, according to the *explanation claim*, people naïvely represent our  
3244 world to be a growing block because they naïvely represent the future to be  
3245 open.

3246 Our aim is not to investigate all these claims in their full generality but  
3247 rather to investigate certain aspects of these claims as they pertain to nomic  
3248 openness.

3249 Consider, first, the narrow vindication claim. In order to evaluate the narrow  
3250 vindication claim, we would need to know the content of our naïve represen-  
3251 tation of future openness. This paper will speak to the issue of whether our  
3252 naïve representation of future openness is partly constituted by our represent-  
3253 ing it to be nomically open. So, it will provide the beginnings of the sort of  
3254 account we would need to determine whether the narrow vindication claim  
3255 (and indeed the vindication claim itself) is true.

3256 Next, consider the reason claim. We investigate whether people take the fact  
3257 that a world is nomically open to be a reason to think that it is a growing block  
3258 world rather than a block universe world. We also investigate a particular view  
3259 about what this reasoning might consist in. According to this view, people  
3260 reason from their ability to deliberate and to act freely to the idea that the  
3261 future is nomically open. They then reason from the nomic openness of the  
3262 future to the idea that future events do not exist because they think that if  
3263 future events did exist “out there in spacetime,” then those events must be  
3264 determined because facts about them already obtain. But in representing  
3265 that future events do not exist and will later come to exist, one represents  
3266 one crucial element of the growing block view. Thus, it might be that by  
3267 representing the world as nomically open, people come to represent it as a  
3268 growing block.

3269 Now, to be clear, we are not endorsing either stage of this reasoning from  
3270 freedom/deliberation to nomic openness or from nomic openness to the non-  
3271 existence of future events (indeed, this last inference is clearly invalid). We  
3272 are merely hypothesising that people (likely tacitly) reason in something like  
3273 this manner, and so they take the presence of nomic openness in a world to  
3274 be a reason to think that the world is a growing block world rather than a  
3275 block universe world. We will call the claim that people reason in this way  
3276 the *deliberative reasoning claim*.

3277 Finally, according to the version of the explanation claim that we investigate  
3278 here, the fact that people naïvely represent the future as nomically open is  
3279 part of what explains why they represent our world to be a growing block.  
3280 Notice that the reason claim and the explanation claim can come apart. It

3281 could be that people naïvely represent our world as a growing block because  
3282 they represent it as nomically open, even though they do not tacitly suppose  
3283 that the latter is a reason to think our world is a growing block (perhaps there  
3284 is a common cause of both representations). Equally, it could be that people  
3285 *do* think that a world being nomically open is a reason to think it is a growing  
3286 block rather than a block universe, but this does not, in fact, explain why  
3287 people think our world is a growing block world (either because they don't  
3288 think it is a growing block, or because they don't think our world is nomically  
3289 open, or because other factors completely swamp this reason and do all the  
3290 explanatory work).

3291 In experiment 1, we seek to determine whether people's naïve representa-  
3292 tion of the future involves nomic openness. We present participants with  
3293 two *nomic vignettes*: one that describes a nomically open world and one that  
3294 describes a nomically closed world. Having seen the two vignettes, partici-  
3295 pants are then asked which world is most like our world (nomically open  
3296 or closed). Our first hypothesis (H1) is that more people will judge that the  
3297 nomically open world is more like our world than the nomically closed world.  
3298 If most people naïvely represent the future as nomically open, then it seems  
3299 reasonable to say that their naïve representation of the future as open consists,  
3300 at least in part, in them representing the future in this matter.

3301 Participants are then presented with two *time vignettes*, one describing a  
3302 growing block world and one describing a block universe world. They are then  
3303 asked which world is most like our world. We predicted (H2) that more people  
3304 would judge that our world is like the growing block world than the block  
3305 universe world. This hypothesis is motivated by previous work on the way that  
3306 people naïvely represent time, including that of Latham, Miller and Norton  
3307 (2021a, 2021b, 2023), and, if vindicated, would replicate these findings.

3308 If the explanation claim is true, then we should find an association between  
3309 people judging that the nomically open world is most like our world and  
3310 judging that the growing block world is most like our world, and between  
3311 people judging that the nomically closed world is most like our world and  
3312 judging that the block universe world is most like our world. This was H3.

3313 In order to investigate the reason claim, we present participants with just  
3314 one of the nomic vignettes. Those who see the nomically open vignette are  
3315 told that Katie is in a world just like that and then asked whether she is more  
3316 likely to be in the growing block or the block universe world. Those who see  
3317 the nomically closed vignette are told that Katie is in a world just like that  
3318 and then asked whether she is more likely to be in the growing block or the

3319 block universe world. If the reason claim is true, then people should judge  
 3320 that if Katie is in a nomically open world, then she is more likely to be in a  
 3321 growing block world as opposed to a block universe world, and if Katie is in a  
 3322 nomically closed world, then they should judge that she is more likely to be  
 3323 in a block universe world as opposed to a growing block world. This was our  
 3324 H4.

3325 Experiment 2 tests the deliberative reason claim. Here, participants are  
 3326 presented with a single vignette that describes an interaction between two  
 3327 characters (George and Helena). George reasons from the fact that our world  
 3328 is deliberatively open to the conclusion that it is nomically open and, from  
 3329 there, to the conclusion that future events do not exist. Helena rejects George's  
 3330 reasoning and explains where she thinks it goes awry. Participants are asked  
 3331 which character is correct. If the deliberative reason claim is true, then we  
 3332 should find that more people will judge that George is correct. This is H5. The  
 3333 final part of this experiment focuses on whether people can see the inferential  
 3334 connection between accepting or rejecting this reasoning. Participants are  
 3335 asked which world (growing block or block universe) the two characters  
 3336 will take *themselves* to be in. We predicted that participants would judge that  
 3337 Helena would take herself to be in a block universe world while George would  
 3338 take himself to be in a growing block world (H6).

3339 We begin, in section 1, by outlining our methodology and results. Then,  
 3340 in section 2, we consider the upshot of those results for understanding our  
 3341 pre-reflective views of the world and the connection between them.

## 3342 1 Methodology and Results

### 3343 1.1 Experiment 1 Methodology

#### 1.1.1 1.1.1 Participants

3345 856 people participated in the study. Participants were recruited and tested  
 3346 online using Amazon Mechanical Turk and compensated \$2 for their time.  
 3347 732 participants had to be excluded from the analyses. That is because they  
 3348 failed to answer all the questions ( $n = 80$ ), failed one of the attentional check  
 3349 questions ( $n = 73$ ), or failed to answer two out of three comprehension ques-  
 3350 tions correctly for the openness vignettes or three out of four comprehension  
 3351 questions correctly for both time vignettes ( $n = 579$ ). The remaining sample  
 3352 was composed of 124 participants (46 female; aged 21–72, mean age 38.98

3353 (SD = 9.95)). Ethics approval for these studies was obtained from the Univer-  
3354 sity of Sydney Human Research Ethics Committee. Informed consent was  
3355 obtained from all participants prior to testing. The survey was conducted  
3356 online using Qualtrics.<sup>8</sup>

## 1.1.2 Materials and Procedure

3358 Participants first see *both* of the following openness vignettes. The first vignette  
3359 describes a world in which the universe is Nominally Open—which we called  
3360 Universe A. The second vignette describes a world in which the universe is  
3361 Nominally Closed—which we called Universe B.

### 3362 NOMICALLY OPEN (UNIVERSE A):

3363 Imagine a universe (Universe A) in which not everything that hap-  
3364 pens is completely caused by whatever happened before it. In Uni-  
3365 verse A, there are multiple different ways the future could go, given  
3366 that the past and present are as they are. Given the past, every event  
3367 *does not have to happen* the way that it does. So, if we “ran” Uni-  
3368 verse A over again from its very first moment, events might unfold  
3369 differently to the way they did unfold.

3370 For example, one day, Katie decided she wanted to have a cup of  
3371 coffee with her breakfast. Like everything else, this decision is not  
3372 completely caused by whatever happened before it. So, if everything  
3373 in the universe was exactly the same up until Katie made her deci-  
3374 sion, it *did not have to happen* that Katie would decide to have a cup  
3375 of coffee.

### 3376 NOMICALLY CLOSED (UNIVERSE B):

3377 Imagine a universe (Universe B) in which everything that happens  
3378 is completely caused by whatever happened before it. In Universe  
3379 B, there are not multiple different ways the future could go, given  
3380 that the past and present are as they are. Given the past, every event  
3381 *has to happen* the way that it does. So, if we “ran” Universe B over  
3382 again from its very first moment, events would unfold exactly the  
3383 same way that they did unfold.

3384 For example, one day, Katie decided she wanted to have a cup of  
3385 coffee with her breakfast. Like everything else, this decision was

---

8 22% of the remaining sample got every comprehension question correct.

3386 completely caused by whatever happened before it. So, if everything  
 3387 in this universe was exactly the same up until Katie made her deci-  
 3388 sion, then it *had to happen* that Katie would decide to have a cup of  
 3389 coffee.

3390 After reading both vignettes, participants responded to three comprehen-  
 3391 sion questions to which they could either respond (a) true or (b) false.

- 3392 1. If we “reran” Universe [A/B] over and over again, we would always get  
 3393 the very same events occurring in the very same order.
- 3394 2. In Universe [A/B], the way things are now could not have been any  
 3395 different from how they are unless the past had been different from  
 3396 how it is.
- 3397 3. In Universe [A/B], there is only one way the future can unfold, given  
 3398 that the past and present are the way they are.

3399 Participants who did not correctly answer two out of three of these questions  
 3400 for each vignette were excluded from the analyses.

3401 Participants are then asked, “Which universe do you think is most like our  
 3402 universe?” and given two options: (a) **UNIVERSE A** or (b) **UNIVERSE B**.

3403 Participants then see both of the following time vignettes. The first vignette  
 3404 describes a universe that is a growing block world—which we called Universe  
 3405 C. The second vignette describes a block universe world—which we called  
 3406 Universe D.

3407 **GROWING BLOCK (UNIVERSE C):**

3408 Imagine a universe (Universe C) where new events—such as the  
 3409 extinction of the dinosaurs, the launching of a ship, or the cutting  
 3410 of a birthday cake—and objects—such as the birth of a baby or the  
 3411 creation of a new car—constantly come into existence. The events  
 3412 and objects that come into existence remain in existence, so the  
 3413 sum total of reality grows as new events and objects come to exist.  
 3414 In this universe, the events and objects that have just come into  
 3415 existence are those that are in the objective present. As new events  
 3416 and objects come into existence, already existing events and objects  
 3417 become part of the past. No future events or objects exist. So, there  
 3418 is a real, objective fact of the matter about which events are present  
 3419 and which are past.

3420 For example, in Universe C, there is the event of Suzy throwing  
3421 the ball at the window and the event of Billy throwing the ball at  
3422 the window. When Suzy throws her ball, Billy is still holding his  
3423 ball; he has yet to throw it. When the event of Suzy's ball hitting  
3424 the window comes into existence, it is in the objective present, and  
3425 the event of Billy's ball hitting the window does not yet exist. It is  
3426 still in the future. When the event of Billy's ball hitting the window  
3427 comes into existence, it is in the objective present, and the event  
3428 of Suzy's ball hitting the window exists in the objective past. So, in  
3429 this universe, first, Suzy throws the ball, and it hits the window;  
3430 then, later, the event of Billy's ball hitting the window comes into  
3431 existence, at which time Suzy's throwing the ball at the window still  
3432 exists but is in the past.

3433 BLOCK UNIVERSE (UNIVERSE D):

3434 Imagine a universe (Universe D) where a single set of events—such  
3435 as the extinction of the dinosaurs, the launching of a ship, or the  
3436 cutting of a birthday cake—and objects—such as the birth of a baby  
3437 or the creation of a new car—exist. All these events are equally real.  
3438 The sum total of reality never grows or shrinks, so the totality of  
3439 events that exist never changes. In this world, past, present, and  
3440 future events all exist. If dinosaurs have ever existed, then dinosaurs  
3441 exist somewhere in the universe. If there will ever exist sentient  
3442 robots, then there exist sentient robots somewhere in the universe.  
3443 In Universe D, other *times* are much like other *places*. Just as in our  
3444 world, Singapore, Sydney, and Seattle all exist, even though they  
3445 do not exist in the same place; in Universe D, dinosaurs and robots  
3446 exist, even though they do not exist at the same time. So, in Universe  
3447 D, every time is present from the perspective of those located in it,  
3448 just as every place is “here” from the perspective of those located in  
3449 it.

3450 For example, in Universe D, there is the event of Suzy throwing  
3451 the ball at the window and the event of Billy throwing the ball at the  
3452 window. When Suzy throws her ball, Billy is still holding his ball;  
3453 he has yet to throw it. In Universe D, the event of Suzy throwing  
3454 her ball and the event of Billy throwing his ball both exist. But they  
3455 do not exist at the same place in space-time: the event of Suzy's ball  
3456 hitting the window is earlier than the event of Billy's ball hitting the

3457 window. So, in Universe D, there is a fact of the matter about which  
 3458 ball hits the window first, namely, Suzy's, and so there is a fact of  
 3459 the matter about in which order the two events occur. But there  
 3460 is no fact about which event *really is* present and which is past or  
 3461 future. The event of Suzy's ball hitting the window is *past* relative to  
 3462 people who are located at the time that Billy's ball hits the window,  
 3463 while the event of Billy's ball hitting the window is *future* relative to  
 3464 people who are located at the time that Suzy's ball hits the window.

3465 After reading both time vignettes, participants responded to four compre-  
 3466 hension questions to which they could respond (a) true or (b) false.

- 3467 1. In Universe [C/D], the past and present exist, but the future does not.
- 3468 2. In Universe [C/D], the past, present, and future exist.
- 3469 3. In Universe [C/D], there is an objective fact as to which events are  
 3470 present.
- 3471 4. In Universe [C/D], events are always past or future relative to other  
 3472 events.

3473 Participants who failed to correctly answer three out of four of these questions  
 3474 for each vignette were excluded from the analyses.

3475 Participants are then asked, "Which universe do you think is most like our  
 3476 universe?" and are given two options: (a) **UNIVERSE C** or (b) **UNIVERSE D**.

3477 Finally, participants then see either the nomically open or nomically closed  
 3478 vignette again, along with both time vignettes, and respond to the following  
 3479 question: "Katie is in a universe just like A/B. Do you think that Katie is more  
 3480 likely to be in Universe C or more likely to be in Universe D?" and are given  
 3481 two options: (a) **UNIVERSE C** or (b) **UNIVERSE D**.

## 1.1.3 Results

3483 Before presenting the statistical analysis, we will start by summarising our  
 3484 main findings. We first hypothesised that (H1) more people would judge that  
 3485 the nomically open world is more like our world than the nomically closed  
 3486 world. This hypothesis was supported. Participants were more likely to judge  
 3487 that our world is more like a nomically open world compared to a nomically  
 3488 closed world. We then hypothesised that (H2) most people would judge that  
 3489 our world is a growing block world rather than a block universe world. This  
 3490 hypothesis was not supported.



3491 Next, we hypothesised (H<sub>3</sub>) that there would be an association between  
3492 people judging that the nomically open world is most like our world and  
3493 judging that the growing block world is most like our world and between  
3494 people judging that the nomically closed world is most like our world and  
3495 judging that the block universe world is most like our world. This hypothesis  
3496 was not supported. While there was a significant association between people's  
3497 judgements about nomic openness and time, the association we found was not  
3498 the one we hypothesised. Instead, there was an association between judging  
3499 that our world is nomically closed and judging it to be a growing block world.  
3500 Participants who judged our world to be nomically open were roughly divided  
3501 in their likelihood to judge our world to be a growing block world or a block  
3502 universe world.

3503 Finally, we hypothesised (H<sub>4</sub>) that participants who are told that a character  
3504 (Katie) is in a nomically open world would be more likely to judge that she is  
3505 in a growing block world than a block universe world (and participants who  
3506 are told that she is in a nomically closed world would be more likely to judge  
3507 that she is in a block universe world than a growing block world). We found  
3508 evidence for this.

3509 Separate one-way chi-square tests were performed to test whether  
3510 (a) most participants judged that the nomically open world was more  
3511 like our world compared to the nomically closed world and whether (b)  
3512 most participants judged that our world is a growing block world rather  
3513 than a block universe world. The results of those tests showed that the  
3514 first hypothesis was vindicated. This means that participants are more  
3515 likely to judge the world as nomically open (76, 61.3%) as opposed to  
3516 being nomically closed (48, 38.7%;  $\chi^2(1, N = 124) = 6.323, p = .012$ ).  
3517 Our hypothesis that participants will judge that our world is more like  
3518 a growing block world (69, 55.9%) as opposed to a block universe world  
3519 (55, 44.4%;  $\chi^2(1, N = 124) = 1.582, p = .209$ ) was not statistically significant,  
3520 indicating that participants are equally likely to judge our world as being  
3521 either a growing block world or a block universe world.

3522 Table 1 below summarises the descriptive data of participants' judgements  
3523 regarding which nomic vignette (**NOMICALLY OPEN**; **NOMICALLY CLOSED**)  
3524 is most like our world and which time vignette (**GROWING BLOCK** world;  
3525 **BLOCK UNIVERSE** world) is most like our world. To test whether there was an  
3526 association between participants who judged our world to be nomically open  
3527 and their judging of our world to be a growing block world, we performed a chi-  
3528 square test of independence. This hypothesis was not supported. Instead, there

3529 was an association between participants judging our world to be nomically  
 3530 closed and judging it to be a growing block world ( $\chi^2(1, N = 124) = 5.449, p =$   
 3531  $.020$ ). Participants who judged our world to be nomically open were divided  
 3532 between judging it to be a growing block world and a block universe world.

Table 1: Participants' judgements of which nomic universe and time vignette are most like the actual world.

<b>World</b>	<b>Growing Block World</b>	<b>Block Universe</b>
<b>Nomically Open</b>	(36) 29.0%	(40) 32.3%
<b>Nomically Closed</b>	(33) 26.6%	(15) 12.1%

3533 Finally, we performed a chi-square test of homogeneity to test whether  
 3534 participants who are told that Katie is in a nomically open world would be  
 3535 more likely to judge that she is in a growing block world (and whether people  
 3536 who are told that she is in a nomically closed world would be more likely to  
 3537 judge that she is in a block universe world). There was a significant association,  
 3538  $\chi^2(1, N = 124) = 6.613, p = .010$ . Participants who were told that Katie was  
 3539 in a nomically open world were more likely to judge that she was also in a  
 3540 growing block world. Meanwhile, participants who were told that Katie was  
 3541 in a nomically closed world were more likely to judge that she was also in a  
 3542 block universe world (see table 2).

Table 2: Participants' judgements of which universe Katie is more likely to be in based on associations between nomic openness and time.

<b>World</b>	<b>Growing Block World</b>	<b>Block Universe</b>
<b>Nomically Open</b>	(38) 65.5%	(20) 34.5%
<b>Nomically Closed</b>	(28) 42.4%	(38) 57.6%

## 1.4.2 Experiment 2 Methodology

### 1.2.4.1 Participants

3545 856 people participated in the study. Participants were recruited and tested  
 3546 online using Amazon Mechanical Turk and compensated \$2 for their time.

3547 732 participants had to be excluded from the analyses. That is because they  
3548 failed to answer all the questions ( $n = 124$ ), failed one of the attentional  
3549 check questions ( $n = 54$ ), or failed to answer three out of four comprehension  
3550 questions correctly for the discussion vignette or failed to answer three out of  
3551 four comprehension questions correctly for the time vignettes ( $n = 554$ ). The  
3552 remaining sample was composed of 124 participants (49 female, 2 trans/non-  
3553 binary; aged 20–78, mean age 36.58 ( $SD = 99.716$ )). Ethics approval for these  
3554 studies was obtained from the University of Sydney Human Research Ethics  
3555 Committee. Informed consent was obtained from all participants prior to  
3556 testing. The survey was conducted online using Qualtrics.<sup>9</sup>

## 1.2.2 Materials and Procedure

3558 In this study, participants first see a single vignette—the nomic discussion  
3559 vignette—in which Helena and George present different views about the  
3560 connection between nomic openness and the existence of the future.

### 3561 NOMIC DISCUSSION:

3562 Helena and George are standing outside a philosophy room having  
3563 a heated discussion about the reasons there are to think that the  
3564 future either exists or does not exist. If the future **does not** exist,  
3565 then future events, such as the existence of a colony on Mars or the  
3566 robot uprising, do not exist, although perhaps one day they will. If  
3567 the future **does** exist, then if there will be a colony on Mars in the  
3568 future, it is true right now that the colony exists out there in the  
3569 universe somewhere. If the future exists, then future events (and  
3570 places) are much like other places here and now. While Helena and  
3571 George are located in Singapore, it's still the case that Sydney and  
3572 London exist; they just don't exist *in Singapore*. In the same way, if  
3573 the future exists, then the colony on Mars exists; it just doesn't exist  
3574 *here and now*.

3575 According to George, one reason to think that the future does  
3576 not exist is that if the future did exist, then there are not multiple  
3577 different ways the future could go, given that the past and present are  
3578 as they are. If the future exists, then given the past and present, every  
3579 future event *has to happen* the way that it does. So if the future exists,  
3580 then if we re-ran the universe over again from its very first moment,

---

9 16% of the remaining sample got every comprehension question correct.

3581 events would unfold exactly the same way. But then Helena cannot  
 3582 be free to *choose* what to eat for breakfast tomorrow, since whatever  
 3583 she eats for breakfast tomorrow, it *had* to be that she would eat that  
 3584 thing.

3585 Helena tells George that he is mistaken. That kind of reasoning,  
 3586 she says, gives us no reason to think that the future does not exist.  
 3587 Just because the event of my (Helena's) eating cereal exists out there  
 3588 in the future, it doesn't mean that my eating cereal was determined  
 3589 by the past and present. It doesn't mean that the future could not  
 3590 have gone some other way. It could be that if we reran the universe  
 3591 over again, then I would instead eat toast instead of cereal for break-  
 3592 fast. The mere fact that the event of my eating cereal is out there  
 3593 in the universe doesn't tell us that that event *had* to be out there.  
 3594 You, George, are located here in this office. But the fact that you are  
 3595 located here doesn't tell us that if the past and present had been the  
 3596 same, you *had* to be located in this office. Perhaps you could have  
 3597 been somewhere different! So, the fact that the event of my eating  
 3598 cereal is out there in the universe does not mean that I *had* to eat  
 3599 cereal. It just means that, in fact, I do eat cereal.

3600 Participants then answered four comprehension questions to which they  
 3601 could answer either (a) true or (b) false.

- 3602 (a) If Helena is right, then if the future exists, it can still be true that there  
 3603 are multiple ways the future could go, given that the past and present  
 3604 are as they are.
- 3605 (b) If George is right, then if the future exists, it can still be true that there  
 3606 are multiple ways the future could go, given that the past and present  
 3607 are as they are.
- 3608 (c) According to Helena, if the event of her eating cereal tomorrow exists,  
 3609 then it could still be that the past and present did not determine that  
 3610 she would decide to eat cereal.
- 3611 (d) According to George, if the event of her eating cereal tomorrow exists,  
 3612 then it must be that the past and present determined that she would  
 3613 decide to eat cereal.

3614 Participants who failed to correctly answer three out of four of these questions  
 3615 were excluded from the analyses.

3616 Participants are then asked, “Which of the two parties, Helena or George,  
3617 do **you** think is right?” and are given two options: (a) George or (b) Helena.

3618 Participants then see both the time vignettes and associated comprehension  
3619 questions (see experiment 1). Participants who failed to correctly answer  
3620 three out of four of these questions for each vignette were excluded from the  
3621 analyses.

3622 Finally, participants then saw the **NOMIC DISCUSSION** vignette again, along  
3623 with both time vignettes. They were then presented with two questions:

- 3624 (1) “Which universe do you think *Helena* will think is most like the universe  
3625 she is in?”  
3626 (2) “Which universe do you think *George* will think is most like the universe  
3627 he is in?”

3628 For each question, they were given two options: (a) **UNIVERSE C** or (b) **UNI-  
3629 VERSE D**.

## 1.2.3 Results

3631 As in experiment 1, we also tested H2 by asking participants which world  
3632 they believed was most like our world (i.e., growing block world or block  
3633 universe world) and predicted that most people would judge that our world is  
3634 a growing block world rather than a block universe world. Again, H2 was not  
3635 supported. People were divided between judging that our world is most like a  
3636 growing block world and a block universe world.

3637 We hypothesised that (H5) if the deliberative reasoning claim is right, then  
3638 most people should judge that George, rather than Helena, is right in the  
3639 **NOMIC DISCUSSION** vignette. This hypothesis was not supported. Instead,  
3640 contrary to our prediction, we found that most participants judged that Helena,  
3641 rather than George, was right.

3642 Finally, we hypothesised that (H6) people will judge that Helena will take  
3643 herself to be in a block universe world and that George will take himself to  
3644 be in a growing block world. This hypothesis was supported.

3645 Separate one-way chi-square tests were performed to test whether (a) most  
3646 participants will judge that our world is more like a growing block world, (b)  
3647 most participants will judge that George was right in the nomic openness  
3648 discussion, (c) most participants will judge that Helena will take herself to be  
3649 in a block universe world, and (d) most participants will judge that George  
3650 will take himself to be in a growing block world. The results of those tests

3651 showed that (a) participants were divided between judging that our world  
 3652 is more like a growing block world (64, 51.6%) and a block universe world  
 3653 (60, 48.4%;  $\chi^2(1, N = 124) = .124, p = .129$ ), which does not support H<sub>2</sub>.  
 3654 Further, (b) contrary to H<sub>5</sub>, more participants judged that Helena (87, 70.2%),  
 3655 rather than George (37, 29.8%), was right in the nomic openness discussion,  
 3656  $\chi^2(1, N = 124) = 20.161, p < .001$ . H<sub>6</sub> was vindicated: most participants  
 3657 (c) judged that Helena would take herself to be in the block universe world  
 3658 (80, 64.5%;  $\chi^2(1, N = 124) = 10.452, p < .001$ ), and that (d) George would  
 3659 take himself to be in the growing block world (80, 64.5%;  $\chi^2(1, N = 124) =$   
 3660  $10.452, p < .001$ ).

## 3662 Discussion

3662 There are several notable aspects of our results. First, as predicted, we found  
 3663 that a majority of people judged our world to be nomically open rather than  
 3664 closed. These results are of interest to those aiming to model our naïve rep-  
 3665 resentation of future openness. Taken in conjunction with previous work in  
 3666 this area, they begin to paint a picture of people's naïve representation of the  
 3667 future.

3668 Hodroj et al. (2023) found that a majority of people (66%) judged our world  
 3669 to be one in which the future is *alethically* open rather than closed. Latham  
 3670 and Miller (2023) report that a majority of people (87%) judged our world to be  
 3671 deliberately open rather than deliberately closed: that is, they judged the  
 3672 future to be one in which what we do in the future is the product of our earlier  
 3673 deliberations, so that had we deliberated differently, we would have made  
 3674 different choices and subsequently done different things. These results, taken  
 3675 together with our current results, suggest that people's naïve representation  
 3676 of the future probably involves at least a combination of representing the  
 3677 future to be deliberately, alethically, and nomically open. It also suggests  
 3678 that it may be deliberative openness that is most important when it comes to  
 3679 capturing people's naïve representation of the open future (something Torre  
 3680 2011 gestures towards).

3681 These results may also suggest that there are several naïve representations  
 3682 of future openness, all or almost all of which include representing the future  
 3683 as deliberately open but only some of which include representing it as  
 3684 nomically and/or alethically open. Perhaps this is not surprising, given the  
 3685 evidence regarding people's naïve representation of time. Baron, Miller and  
 3686 Tallant (2022) cite a range of experiments that they take jointly to show that

3687 there is no single, shared, naïve representation of time. What is true of time  
3688 might also be true of naïve representations of the open future.

3689 Our results also have implications for the narrow vindication claim. Ac-  
3690 cording to that claim, recall, the growing block theory vindicates our naïve  
3691 representation of the future as open, and the block universe theory does not.  
3692 There is some support for this claim, given the results of this study, alongside  
3693 those of Hodroj et al. (2023) and Latham et al., despite the fact that these  
3694 studies jointly suggest that *most* aspects of our naïve representation of future  
3695 openness (and the most important of these) are consistent with our world  
3696 being a block universe world.

3697 The study by Latham et al. suggests that a vast majority of people have naïve  
3698 representations of the future according to which the future is deliberately  
3699 open. But the presence of deliberative openness is clearly consistent with our  
3700 world being either a block universe or a growing block world. So, arguably,  
3701 the most powerful aspect of our naïve representation of the future is one that  
3702 can be vindicated by either view of time.

3703 The current study found that a majority of people represent the future as  
3704 nomically open, not closed. But, again, the future being nomically open is  
3705 consistent with our world being either a block universe or a growing block.  
3706 So, either view can vindicate this aspect of our naïve representation.

3707 The only good news for the growing block theorist lies in the Hodroj et  
3708 al. (2023) study, which found that a majority of people represent the future  
3709 as alethically open. On standard (i.e., nonbranching) versions of the block  
3710 universe, the future is not alethically open, while on standard versions of the  
3711 growing block theory, it is. So, the growing block theory does vindicate *this*  
3712 aspect of openness, while the block universe view does not.

3713 Still, it's worth bearing in mind that according to the study by Hodroj et  
3714 al. (2023), ~34% of people did not judge the future to be alethically open. So,  
3715 it may be that a substantial minority of people have a naïve representation  
3716 of the future that is equally vindicated by both the growing block and block  
3717 universe theories. And, of course, even if the narrow vindication claim is true,  
3718 it remains open to dispute whether it gives us much, if any, reason to prefer  
3719 the growing block view to the block universe view. Still, these studies suggest  
3720 that, insofar as growing block theorists want to try and argue for their view  
3721 via something like the (narrow) vindication claim, they might do well to focus  
3722 more on alethic openness than other forms of openness.

3723 Moving on, we did not find that a majority of people represent our world as a  
3724 growing block rather than a block universe. Instead, across both experiments,

3725 people were evenly split between the two models. This should, perhaps, not  
3726 be such a surprise. Latham, Miller and Norton (2021a) found that across  
3727 two experiments, ~70% of people judged our world to be dynamical (either  
3728 growing block, moving spotlight, or presentist), and of those, between ~35%  
3729 and ~50% judged it to be a growing block. Even though in these studies only  
3730 ~25% and ~35% of all people judged our world to be most like a growing block  
3731 world, we expected that, given a forced choice between a growing block and a  
3732 block universe world, most people would judge it to be *more like* a growing  
3733 block world than a block universe world, given that most people judge our  
3734 world to be temporally dynamical.

3735 Our results suggest that although people are drawn to dynamical theories  
3736 of time, their naive representation of time might be less *strongly* dynamical  
3737 than has otherwise been thought. This might explain why, given that the block  
3738 universe and growing block views are very similar in a number of ways, when  
3739 given a forced choice between the two, people tended to be roughly evenly  
3740 divided in which world they thought was most like ours.

3741 This brings us to the explanation and reason claims. Our results here are  
3742 both startling and puzzling. Consider, first, the explanation claim. Our hy-  
3743 pothesis here (H<sub>3</sub>) was not vindicated. While we did find an association, it  
3744 was the opposite of the one we predicted. We found an association between  
3745 judging a world to be nomically *closed* and judging it to be a growing block  
3746 world. Amongst people who judged our world to be nomically open, people  
3747 were evenly split between judging it to be a growing block or a block universe.  
3748 While the latter absence of an association is not such a surprise (given that,  
3749 *in fact*, nomically open words are no more likely to be growing block worlds  
3750 as opposed to block universe worlds, it is perhaps heartening to see people's  
3751 judgements in this regard), the presence of the converse association is very  
3752 puzzling. It's hard to see why people who judge the future to be nomically  
3753 *closed* would tend to judge it to be a growing block. The best we can come  
3754 up with is that perhaps some people think that the laws of nature "push" the  
3755 world along and cause it to grow, and they imagine this growth process must  
3756 be deterministic (else the world would not know what to grow into). If this is  
3757 the reason why (some) people judge our world to be nomically closed, then  
3758 we would expect those people to judge that our world is a growing block. All  
3759 we can really say is that further investigation of the association here would be  
3760 useful.

3761 Certainly, though, the lack of any association between people judging our  
3762 world to be nomically open and judging it to be a growing block world suggests



3763 that it is unlikely that the fact that people naïvely represent the future as  
3764 nomically open is what even partially explains why they represent it to be  
3765 a growing block. This finding is interesting, given our results regarding the  
3766 reason claim. Our hypothesis in this regard was vindicated: participants judged  
3767 that Katie was more likely to be in a growing block world than a block universe  
3768 world if she was in a nomically open world and to be in a block universe  
3769 rather than a growing block world if she was in a nomically closed world.  
3770 Thus, people do seem to think that the fact that a world is nomically open is a  
3771 reason to think it is a growing block world rather than a block universe world.  
3772 The reason claim seems to be vindicated.

3773 The vindication of the reason claim does suggest that there is some *sense*  
3774 in which the growing block view of time better accords with our naïve repre-  
3775 sentation of the future as nomically open. It accords in at least this sense: if  
3776 the only thing someone knows about a world is that it is nomically open, they  
3777 will think it more likely that the world is a growing block rather than a block  
3778 universe world. So, there is some important connection between people's  
3779 naïve representation of the future and their naïve representation of time. The  
3780 former, we might say, *predisposes* them to thinking that our world is a growing  
3781 block world, since if all they know about our world is that it is nomically open,  
3782 people will tend to judge that it is a growing block world.

3783 But of course, this is not all that people know about our world, and presuma-  
3784 bly, this explains why we found no association between people judging that  
3785 our world is nomically open and that it's a growing block world. One thought  
3786 about what might be going on here is that contemporary scientific knowledge  
3787 is pushing people who judge that our world is nomically open to judge that it  
3788 is a block universe world rather than a growing block world. If so, that could  
3789 tend to eliminate the predicted association. But, first, we know from previous  
3790 research by Latham, Miller and Norton (2021a) that levels of education and  
3791 levels of scientific knowledge, especially in physics, have no effect on people's  
3792 judgements about which view of time they think is true of our world. Second,  
3793 in this study, we found that ~50% of people judged our world to be a growing  
3794 block. So, it seems unlikely that this explains why we found no association.

3795 Another possibility is that the reason at least some people judge our world  
3796 to be nomically open is that they are aware of quantum mechanics rather than  
3797 on the basis of their judgement on their naïve representation of the future. If  
3798 so, it may be that those who *naïvely* represent the future as nomically open  
3799 *are* more inclined to represent it as a growing block, but that many of those  
3800 who represent the future as nomically open are employing a scientifically

3801 informed representation of the future, and, perhaps, those people also tend  
3802 to represent the world as a block universe. If so, that could eliminate the  
3803 association. It would be useful to do follow-up work here that attempts to  
3804 determine to what extent people's representation of the future as nomically  
3805 open is naïve, as opposed to scientifically informed.

3806 What we can say, though, is that at best, people are predisposed to represent  
3807 our world to be a growing block in virtue of representing it to be nomically  
3808 open, but that as a matter of fact, what explains why people represent the  
3809 world to be a growing block is not that they represent it to be nomically open.  
3810 This is further suggested by the results of our second experiment, in which  
3811 only ~30% of people judged that George's reasoning was correct. Most people,  
3812 then, do not endorse the deliberative reasoning claim we investigated.

3813 In all, then, we think there is little evidence for the idea that part of what  
3814 explains why people naïvely represent our world as a growing block is that they  
3815 naïvely represent the future as nomically open. This will be of interest to A- and  
3816 B-theorists alike. B-theorists have recently resisted what has become known as  
3817 the argument from temporal phenomenology (Baron et al. 2015)—according  
3818 to which we have reason to think our world is temporally dynamical because  
3819 this is how it seems to us to be in perceptual experience—by denying that  
3820 it does seem this way to us in experience (Hoerl 2014; Prosser 2016; Deng  
3821 2013, 2018; Bardon 2013; Miller, Holcombe and Latham 2020; Miller 2019,  
3822 2023; Latham, Miller and Norton 2020). Such views have often been deemed  
3823 deflationist.

3824 We know, however, that people naïvely represent our world as temporally  
3825 dynamical (Latham, Miller and Norton 2021a, 2021b, 2023). If, as deflationists  
3826 suppose, it does not seem to us, in experience, as though time is dynamical  
3827 (and there is some suggestion from Latham, Miller and Norton 2020 that this  
3828 might be right), then the question arises as to why we naïvely represent it that  
3829 way. Deflationists, it seems, owe us some kind of explanation here.

3830 One possibility, alluded to by Prosser (2016), is that part of what explains  
3831 why we represent time as dynamical is that we represent the future as open.  
3832 This study had the potential to show that part of what explains why we  
3833 represent time as dynamical (by representing it as a growing block) is that we  
3834 represent it as nomically open. Unfortunately for deflationists, we found no  
3835 evidence of this.

3836 Having said that, Prosser's suggestion is rather different from the one we  
3837 investigated here. He hypothesises that because people represent the future  
3838 as being objectively open (as opposed to merely perspectively or subjectively

3839 open), and because we represent that this openness moves (as what was once  
3840 open becomes closed and part of the past), we must represent that there  
3841 is a privileged and moving moment in time that is the border between the  
3842 closed past and the open future. Further work, taking up the specific details of  
3843 Prosser's view, would be welcome, given that we found no evidence in favour  
3844 of the hypotheses we tested in this regard.

3845 In all, we think that there is much more that can be learned about both our  
3846 naïve representation of the open future and the ways in which this representa-  
3847 tion connects to our naïve representation of time. That work can shed light on  
3848 the best way to model future openness (insofar as that modelling is attempting  
3849 to capture some naïve representation of the future) and on whether what  
3850 explains our open future practices also vindicates our naïve representation of  
3851 the open future. It can, we hope, also shed light on the connection between our  
3852 naïve representation of the future and of time and, hence, on extant debates  
3853 in the philosophy of time.\*

3854 Batoul Hodroj  
3855  0000-0002-5227-392X  
3856 University of Queensland  
3857 batoul.hodroj@sydney.edu.au

3858 Andrew James Latham  
3859  0000-0001-8639-8053  
3860 Aarhus University  
3861 andrew.latham@cas.au.dk

3862 Kristie Miller  
3863  0000-0002-5092-8419  
3864 University of Sydney  
3865 kristie.miller@sydney.edu.au

## 3866 References

- 3867 BARDON, Adrian. 2013. *A Brief History of the Philosophy of Time*. Oxford: Oxford  
3868 University Press, doi:10.1093/acprof:oso/9780199976454.001.0001.
- 3869 BARON, Sam, CUSBERT, John, FARR, Matt, KON, Maria and MILLER, Kristie. 2015.  
3870 "Temporal Experience, Temporal Passage and the Cognitive Sciences." *Philosophy*  
3871 *Compass* 10(8): 560–571, doi:10.1111/phc3.12244.

---

\* THANKS

- 3872 BARON, Sam, MILLER, Kristie and TALLANT, Jonathan. 2022. *Out of Time: A Philo-*  
 3873 *sophical Study of Timelessness*. Oxford: Oxford University Press, doi:10.1093/oso/  
 3874 9780192864888.001.0001.
- 3875 BELNAP, Nuel D., Jr. 1992. "Branching Space-Time." *Synthese* 92(3): 385–434, doi:10.1  
 3876 007/bf00414289.
- 3877 —. 2005. "Agents and Agency in Branching Space-Times." in *Logic, Thought and Action*,  
 3878 edited by Daniel VANDERVEKEN, pp. 291–313. Logic, Epistemology, and the Unity  
 3879 of Science n. 2. Dordrecht: Springer, doi:10.1007/1-4020-3167-X\_14.
- 3880 BRADDON-MITCHELL, David. 2004. "How Do We Know it is Now Now?" *Analysis*  
 3881 64(3): 199–203, doi:10.1111/j.0003-2638.2004.00485.x.
- 3882 BRIGGS, Rachael A. and FORBES, Graeme A. 2012. "The Real Truth About the Unreal  
 3883 Future." in *Oxford Studies in Metaphysics*, volume VII, edited by Karen BENNETT  
 3884 and Dean W. ZIMMERMAN, pp. 257–304. New York: Oxford University Press,  
 3885 doi:10.1093/acprof:oso/9780199659081.003.0009.
- 3886 BROAD, Charlie Dunbar. 1923. *Scientific Thought*. International Library of Psychology,  
 3887 Philosophy and Scientific Method. London: Kegan Paul, Trench, Trübner & Co.
- 3888 —. 1938. *Examination of McTaggart's Philosophy. Volume II, Part I*. Cambridge: Cam-  
 3889 bridge University Press.
- 3890 CALLENDER, Craig. 2017. *What Makes Time Special?* Oxford: Oxford University Press,  
 3891 doi:10.1093/oso/9780198797302.001.0001.
- 3892 CORREIA, Fabrice and ROSENKRANZ, Sven. 2018. *Nothing to Come: A Defence of the*  
 3893 *Growing Block Theory of Time*. Synthese Library n. 395. Dordrecht: Springer, doi:10  
 3894 .1007/978-3-319-78704-6.
- 3895 DENG, Natalja. 2013. "On Explaining Why Time Seems to Pass." *The Southern Journal*  
 3896 *of Philosophy* 51(3): 367–382, doi:10.1111/sjp.12033.
- 3897 —. 2018. "On 'Experiencing Time': A Response to Simon Prosser [on Prosser (2016)]."  
 3898 *Inquiry* 61(3): 281–301, doi:10.1080/0020174x.2017.1322674.
- 3899 FORBES, Graeme A. 2016. "The Growing Block's Past Problems." *Philosophical Studies*  
 3900 173(3): 699–709, doi:10.1007/s11098-015-0514-1.
- 3901 FORREST, Peter. 2004. "The Real but Dead Past: A Reply to Braddon-Mitchell (2004)."  
 3902 *Analysis* 64(4): 358–362, doi:10.1111/j.0003-2638.2004.00510.x.
- 3903 GRANDJEAN, Vincent. 2021. "How Is the Asymmetry Between the Open Future and the  
 3904 Fixed Past to Be Characterized?" *Synthese* 198(3): 1863–1886, doi:10.1007/s11229-  
 3905 019-02164-2.
- 3906 —. 2022. *The Asymmetric Nature of Time: Accounting for the Open Future and the Fixed*  
 3907 *Past*. Synthese Library n. 468. Dordrecht: Springer, doi:10.1007/978-3-031-09763-8.
- 3908 HODROJ, Batoul, LATHAM, Andrew James, LEE-TORY, Jordan and MILLER, Kristie.  
 3909 2023. "Alethic Openness and the Growing Block Theory of Time." *The Philosophical*  
 3910 *Quarterly* 73(2): 532–556, doi:10.1093/pq/pqac062.
- 3911 HOERL, Christoph. 2014. "Do We (Seem To) Perceive Passage?" *Philosophical Explor-*  
 3912 *ations* 17(2): 188–202, doi:10.1080/13869795.2013.852615.

- 3913 ISMAEL, Jenann. 2012. "Decision and the Open Future." in *The Future of the Philosophy*  
3914 *of Time*, edited by Adrian BARDON, pp. 149–168. Routledge Studies in Metaphysics  
3915 n. 4. London: Routledge, doi:10.4324/9780203338315.
- 3916 LATHAM, Andrew James and MILLER, Kristie. 2023. "Why Do People Represent Time  
3917 as Dynamical? An Investigation of Temporal Dynamism and the Open Future."  
3918 *Philosophical Studies* 180(5): 1717–1742, doi:10.1007/s11098-023-01940-8.
- 3919 LATHAM, Andrew James, MILLER, Kristie and NORTON, James. 2020. "An Empirical  
3920 Investigation of Purported Passage Phenomenology." *The Journal of Philosophy*  
3921 117(7): 353–386, doi:10.5840/jphil2020117722.
- 3922 —. 2021a. "Is Our Naïve Theory of Time Dynamical?" *Synthese* 198(5): 4251–4271,  
3923 doi:10.1007/s11229-019-02340-4.
- 3924 —. 2021b. "An Empirical Investigation of the Role of Direction in Our Concept of  
3925 Time." *Acta Analytica* 36(1): 25–47, doi:10.1007/s12136-020-00435-z.
- 3926 —. 2023. "Do the Folk Represent Time as Essentially Dynamical?" *Inquiry* 66(10):  
3927 1882–1913, doi:10.1080/0020174X.2020.1827027.
- 3928 LEWIS, David. 1979. "Counterfactual Dependence and Time's Arrow." *Noûs* 13(4):  
3929 455–476. Reprinted, with a postscript (Lewis 1986b), in Lewis (1986a, 32–51),  
3930 doi:10.2307/2215339.
- 3931 —. 1986a. *Philosophical Papers, Volume 2*. Oxford: Oxford University Press, doi:10.109  
3932 3/0195036468.001.0001.
- 3933 —. 1986b. "Postscript to Lewis (1979)." in *Philosophical Papers, Volume 2*, pp. 52–66.  
3934 Oxford: Oxford University Press, doi:10.1093/0195036468.003.0002.
- 3935 MACFARLANE, John. 2003. "Future Contingents and Relative Truth." *The Philosophical*  
3936 *Quarterly* 53(212): 321–336, doi:10.1111/1467-9213.00315.
- 3937 —. 2008. "Truth in the Garden of Forking Paths." in *Relative Truth*, edited by Manuel  
3938 GARCÍA-CARPINTERO and Max KÖLBEL, pp. 81–102. Oxford: Oxford University  
3939 Press, doi:10.1093/acprof:oso/9780199234950.003.0004.
- 3940 MARKOSIAN, Ned. 1995. "The Open Past." *Philosophical Studies* 79(1): 95–105, doi:10  
3941 .1007/bf00989786.
- 3942 MCCALL, Storrs. 1994. *A Model of the Universe: Space-Time, Probability, and Decision*.  
3943 Clarendon Library of Logic and Philosophy. Oxford: Oxford University Press,  
3944 doi:10.1093/acprof:oso/9780198236221.001.0001.
- 3945 MILLER, Kristie. 2019. "Does It Really Seem to Us as Though Time Passes?" in *The*  
3946 *Illusions of Time: Philosophical and Psychological Essays on Timing and Time*  
3947 *Perception*, edited by Valterri ARSTILA, Adrian BARDON, Sean Enda POWER, and  
3948 Argiro VATAKIS, pp. 17–34. London: Palgrave Macmillan, doi:10.1007/978-3-030-  
3949 22048-8\_2.
- 3950 —. 2023. "Against Passage Illusionism." *Ergo* 9(45): 1233–1263, doi:10.3998/ergo.2914.
- 3951 MILLER, Kristie, HOLCOMBE, Alex and LATHAM, Andrew James. 2020. "Temporal  
3952 Phenomenology: Phenomenological Illusion Versus Cognitive Error." *Synthese*  
3953 197(2): 751–771, doi:10.1007/s11229-018-1730-y.

- 3954 PROSSER, Simon. 2016. *Experiencing Time*. Oxford: Oxford University Press, doi:10.1  
3955 093/acprof:oso/9780198748946.001.0001.
- 3956 TOOLEY, Michael. 1997. *Time, Tense, and Causation*. Oxford: Oxford University Press,  
3957 doi:10.1093/0198250746.001.0001.
- 3958 TORRE, Stephan. 2011. "The Open Future." *Philosophy Compass* 6(5): 360–373, doi:10  
3959 .1111/j.1747-9991.2011.00395.x.
- 3960 WILLIAMS, J. Robert G. 2008. "Aristotelian Indeterminacy and the Open Future."  
3961 Unpublished manuscript, <https://philpapers.org/rec/WILAIA-4>.

PROOF

3962

Published by *Philosophie.ch*

3963

Verein philosophie.ch

3964

Fabrikgässli 1

3965

2502 Biel/Bienne

3966

Switzerland

3967

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

3968

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

3969

ISSN 0012-2017

3970

ISBN 1234-5678

3971

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

3972

3973

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

3974

3975

Abstracting and Indexing Services

3976

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

3977

Contents, Current Mathematical Publications, Dietrich's Index

3978

Philosophicus, IBZ — Internationale Bibliographie der Geistes- und

3979

Sozialwissenschaftlichen Zeitschriftenliteratur, Internationale Bibliographie

3980

der Rezensionen Geistes- und Sozialwissenschaftlicher Literatur, Linguistics

3981

and Language Behavior Abstracts, Mathematical Reviews, MathSciNet,

3982

Periodicals Contents Index, Philosopher's Index, Repertoire Bibliographique

3983

de la Philosophie, Russian Academy of Sciences Bibliographies.

# Contents

3984

3985	JP SMIT & FILIP BUEKENS, <i>Is Somaliland a Country?: An Essay on Institutional Objects in the Social Sciences</i> . . . . .	1
3986		
3987	LI ZHANG & LEON HORSTEN, <i>The Minimalist Theory of Truth and the Generalisation Problem</i> . . . . .	23
3988		
3989	FR. JAMES DOMINIC ROONEY, OP, <i>The Problem of Thomistic Parts</i> . . . . .	45
3990	WOLFGANG SPOHN, <i>A Generalization of the Reflection Principle</i> . . . . .	73
3991	BATOUL HODROJ, ANDREW JAMES LATHAM & KRISTIE MILLER, <i>Our Naïve Representation of Time and of the Open Future</i> . . . . .	97
3992		

PROOF