

1 First-person interventions and the meta-problem of consciousness

2 Colin Klein (The Australian National University)

3 Andrew B. Barron (Macquarie University)

4
5
6 **Abstract**

7 Chalmers' (2018) meta-problem of consciousness emphasizes unexpected common ground
8 between otherwise incompatible positions. We argue that the materialist should welcome
9 discussion of the meta-problem. We suggest that the core of the meta-problem is the seeming
10 arbitrariness of subjective experience. This has an unexpected resolution when one moves to
11 an interventionist account of scientific explanation: the same interventions that resolve the
12 hard problem should also resolve the meta-problem.

13
14 *1 The seeming arbitrariness of phenomenal properties*

15 In a vivid thought experiment, Herbert Feigl (1958) imagined an 'autocerebroscope' that
16 would allow us to view, in real time, the brain processes responsible for particular
17 phenomenal states. Much of the subsequent debate over consciousness can be rephrased in
18 terms of how satisfying the autocerebroscope would be, and what that answer would imply.

19
20 Vivid visual evidence often makes scientific explanations compelling. Yet many (if not Feigl)
21 have the intuition that an autocerebroscope would still leave consciousness something of a
22 mystery. Explaining the intuition that there would be such a residue forms the core of what
23 Chalmers (2018) calls the *meta-problem of consciousness*.

24
25 We are realists about consciousness, like Chalmers. We are also materialists and naturalists
26 (Barron and Klein 2016), unlike Chalmers. At a first pass, we think that the meta-problem
27 arises due to limited access to complex brain states, along with a false belief that we have
28 complete access (Lashley 1960; Hilbert 1987, Armstrong 1999; Pettit 2003). Because our
29 awareness extends only to the contents of consciousness and not to the complex mechanisms
30 that support it, it is natural to *feel* that the hard problem is hard. That's where the meta-
31 problem comes in, and what supports the first-person aspect of the hard problem.

33 Yet as Chalmers rightly notes (2018; 23), lack of access can only go so far as an explanation.
34 There are many things in the world to which we lack complete access. When science reveals
35 their essential natures we feel satisfied, not puzzled. Consciousness is not like this.

36

37 In one sense, the lack-of-access model does predict recalcitrance of our intuitions. The core
38 of the hard and meta- problems is a certain feeling of *arbitrariness* about subjective
39 experiences. Why should hearing middle C feel like *that*, rather than something else? A
40 whole family of thought experiments emphasize the degree to which it seems that
41 phenomenal properties could be swapped, eliminated, or otherwise messed with, all with no
42 deep consequence to how we get around in the world (Jackson 1982; Chalmers 1996).

43

44 These thought experiments highlight that it doesn't seem very important *which* subjective
45 experiences we have: were they consistent over an individual's life, phenomenal blue and
46 green (say) could be swapped with no effect. That means that generalizations connecting one
47 brain state to a feeling of green and another to a feeling of blue strike us as arbitrary: there's
48 no reason (it seems) why they couldn't be otherwise.

49

50 Arbitrary relationships are anathema to scientific explanation. So even if we have a good
51 solution to the hard problem, it won't *seem* like a good solution. That is, we think, the real
52 bite of the meta-problem.

53

54

55 2) *Explanation and intervention*

56 The meta-problem, as we've sketched it, is deeply bound up in issues of scientific
57 explanation. We think that the way out requires recognizing that psychophysical relationships
58 are contingent, and that the right kinds of intervention could make them otherwise. This may
59 seem odd to some philosophers, but it follows naturally from developments around scientific
60 explanation.¹

61 The past two decades of philosophy of science have emphasized the importance of direct
62 intervention for explanation (Woodward 2003). Most scientific disciplines care about
63 intervention, and the hunt for control variables is key (Campbell 2007). Even if we cannot

¹ For a more developed story, see Klein and Barron (ms).

64 perform interventions due to our contingent limitations, the sort of information that explains
65 is precisely the sort of information that would permit intervention if we had sufficient power.

66

67 Interventionism is a departure from earlier approaches, especially the deductive-nomological
68 (DN) theory of explanation favored by the positivists (Salmon 1989). The DN theory says
69 that explanation comes from exceptionless covering laws. The DN model arguably lingers on
70 in the background assumptions of consciousness studies, particularly in the search for Neural
71 Correlates of Consciousness.

72

73 The weaknesses of the DN model of explanation are well-known (Salmon 1989). We won't
74 recap them. Instead, we highlight two distinctive features of the interventionist model. First,
75 interventionism provides for explanation without appeal to exceptionless universal laws: the
76 invariant generalizations it appeals to hold only over a certain range of conditions
77 (Woodward 2003). Second, interventionism is fundamentally contrastive. One doesn't
78 explain why X holds *tout court*, but only why X holds *rather than* not-X, or *rather than* {Y
79 or Z or...}. Different contrast classes thus give different explanations (Hitchcock 1996).

80

81 Interventionism has been overlooked in consciousness studies, we suspect, because
82 interventions have mostly been studied in the context of explanations of event-types. In the
83 case of consciousness, this would consist of interventions that change some brain activity *B*
84 to *B**, thereby changing phenomenal state *P* to *P**. In that capacity, there is probably little of
85 philosophical interest to consciousness studies. Knowing that stimulating *this* part of the
86 brain gives rise to the taste of a ham sandwich might be some evidence against the crudest
87 forms of substance dualism, but all parties in the current debate can accommodate it.

88

89 However, while less emphasized, it seems clear that there should also be interventionist
90 explanations of invariant generalizations themselves. Consider Woodward's (2003; 12-13)
91 example of a block sliding down a ramp. Woodward presents the standard derivation of the
92 block's acceleration as an explanation in terms of how acceleration depends on friction. But
93 one might equally well treat that demonstration as an explanation *of the invariant*
94 *relationship itself*. Given the explanation, we can *also* show how this generalization would
95 vary given changes in (say) wind resistance on the block, or if the block ceased to be a solid
96 object and acted as a viscous liquid. That is, the standard derivation is not explanatory simply
97 because it is a *derivation* (this is the lesson of attacks on the DN model). Rather, it is

98 explanatory because it shows why the generalization is one way *rather than other ways it*
99 *could be.*

100

101 This explanatory strategy is, note, a completely natural extension of the two interventionist
102 principles above. First, generalizations are not exceptionless and universal, so it makes sense
103 to ask why a generalization is one way rather another. (Conversely, fundamental physical
104 laws might just be brute facts; they can't be explained because they couldn't be made
105 otherwise.) Second, because explanation is fundamentally contrastive, we can equally well
106 ask about contrasts between generalizations as we can between event-types.

107

108 We think this logic follows over to the explanation of consciousness. Consider some brain
109 state *B* that's responsible for a particular experience *P* of pain. At a minimum, the
110 interventionist says, we should be able to change *P* by changing *B*.² *But we should also be*
111 *able to change the relationship between B and P.* That is, if the *B-to-P* relationship is not just
112 a brute law of nature – and the experimentalist is committed to the idea that it isn't – then it
113 should *also* be a target of intervention.

114

115 These are defeasible presumptions. It could turn out that there is no interesting way in which
116 to manipulate the *B-to-P* relationship. Or it could be that the only interventions possible are
117 crude ones that merely remove all consciousness. These would still be interventions on the *B-*
118 *to-P* relationship, but only in the slightly degenerate sense that pulling the plug on a radio
119 makes a difference to the relationship between the volume knob and the loudness of the
120 music (Woodward 2010). If so, non-materialist stories would gain traction. Conversely, the
121 more systematic and specific the control we gain over the *B-to-P* relationship, the more we
122 can explain.

123

124 *3 The Upgrade*

125

126 So here's the first upshot: to solve the hard problem, we must learn to manipulate the *B-to-P*
127 relationship. Yet that alone won't be quite enough to fix the meta-problem. So long as I
128 cannot *experience* how the *B-to-P* relationship can be varied, it will continue to seem to me

² Note that there might be multiple regions which affect *P*, and multiple ways in which *P* can be affected (Klein 2017). We elide this complication here, though we think it is an important feature of the contrastive aspect of explanation.

129 that this relationship is as arbitrary as fundamental law of physics. We don't escape the
130 vicious circle so easy; the meta-problem still has bite. Put another way, the hard problem
131 may not be something we can be *reasoned* out of, no matter how compelling the evidence.

132

133 Yet so phrased, it is clear that the hard problem stems from limitations in our capabilities,
134 rather than from the nature of consciousness itself. Those limitations could change. For this
135 feeling of arbitrariness is a matter not of lack of knowledge but of lack of *control*: we don't
136 have the right sort of influence over our own brains, and thereby over our phenomenal states.
137 But that is a contingent matter.

138

139 Suppose we upgraded the autocerebroscope so that we could both observe *and change* our
140 own neural activity. So for any *B-to-P* relationship, we could alter it smoothly and
141 seamlessly, in real time, and experience the results. We can imagine that we've become so
142 fluent at this that the move from textbook description to altered experience is as smooth as
143 the move from score to notes seems for a skilled pianist.

144

145 Were we to gain such control, we postulate, we would not be particularly impressed by the
146 hard problem of consciousness. Having ceased to be a mere bystander in the production of
147 our own experiences, they would cease to seem arbitrary. The psychophysical generalizations
148 with which we work, having lost their arbitrary tinge, would appear to us as they are: as
149 explanations of why and how consciousness arises from the brain.

150

151 Present technology doesn't allow this kind of tight coupling between specific brain
152 intervention and specific first person experience, of course. The available methods for
153 intervening on the brain tend to be crude, relatively slow, and relatively dangerous. So the
154 upgraded autocerebroscope will remain a philosopher's fantasy for the near future.

155

156 Yet the tight coupling envisioned is probably unnecessary: the ability to observe and change
157 psychophysical links, even at a relatively coarse grain, might be all we need to tackle the hard
158 problem. For the key, note, is that self-intervention solves two problems at once. It solves an
159 objective explanatory problem, by showing the sort of interventions that affect consciousness.
160 And it solves the subjective problem of understanding why the objective interventions are not
161 simply arbitrary, because it puts us in control. Thus it would solve the meta-problem and the

162 hard problem in one go—emphasizing, as Chalmers does, that the two problems are also
163 intimately linked.

164

165 *4) Conclusion*

166 Where does that leave us with respect to both the hard and the meta-problem of
167 consciousness? A brief recap. We've argued that the core of the hard problem is a feeling of
168 arbitrariness about the mind-body relationship, born out of a certain lack of access to
169 anything that would manipulate those grounds. The very same architectural features that give
170 rise to hard problem also explain why we have the intuitions which constitute the meta-
171 problem. This is (we hope) a characterization of the hard problem that is compatible with
172 Chalmers' topic-neutral criterion (2018; 15ff). It is, at least, the sort of thing that many
173 different theorists can subscribe to: the dualist thinks that the lack of access is a deep
174 metaphysical feature of the world, the materialist realist can tell a variety of stories about
175 why it might be, and the illusionist thinks that the seeming arbitrariness is part and parcel
176 with the overall illusory features of consciousness itself.

177

178 We assume, again, that this story is compatible with a realist position about consciousness.
179 We take ourselves to be realists about consciousness in the same sense that Bohr was a realist
180 about atoms. That is, we think conscious states exist, but we're mostly wrong about their
181 properties and maybe entirely wrong about their essential properties. The meta-problem
182 stems from such an error. On some ways of cashing this out, this might sound like a form of
183 illusionism.³ Yet we agree with Chalmers that illusionism sounds odd---much more so than
184 the prospect of surprises about the nature of the phenomenal.

185

186 We think that there is a useful analogy between the experimental program we suggest and
187 other historical advances in science. The concept of *life* was once as fraught as that of
188 consciousness. It seemed to many (including many philosophers) that there was a
189 fundamental mismatch between the properties of living and nonliving matter, severe enough
190 that the latter could never give rise to the former on its own.

191

³ It's often surprising to non-philosophers that questions that appear to be about consciousness often turn on debates about meaning and reference. Unsurprisingly, we are the sorts of naturalists who also like causal theories of reference.

192 We now know that this impression was a mistake. It stemmed from inadequate concepts of
193 both life and matter. Yet the history by which this mistake was corrected is primarily one of
194 *experimentation* rather than mere observation or rational reflection. Wöhler’s synthesis of
195 urea and Bernard’s dramatic manipulations of homeostatic mechanisms were important
196 demonstrations precisely because they showed how to manipulate what looked like brute
197 facts. Indeed, demonstrations like Wöhler’s were important not because they provided
198 decisive evidence against vitalism (which lingered on for a long time afterwards) but because
199 they gave evidence that the search for non-vitalist explanations might even be possible
200 (Ramberg 2000).

201

202 The position we sketch also has some antecedents in analytic philosophy, usually in oblique
203 discussions of the effect of psychedelic drugs (Langlitz 2016). The closest parallel might be
204 with claims by Thomas Metzinger, for example that:

205 ... scientific research programs on consciousness and its neurofunctional correlates
206 *could* be greatly optimized if researchers were well traveled in phenomenal state
207 space, if they were cultivated in terms of the richness of their own inner experience as
208 well. But not because this would give them a mysterious kind of first- person
209 “data”—more likely, because it would thoroughly shatter their folk-
210 phenomenological intuitions and endow them with completely new *theoretical*
211 intuitions. What is right is that first-person approaches possess an enormous *heuristic*
212 potential, and that we are currently far from realizing it. (Metzinger 2006; 2-3)

213

214 We think there is much to endorse in this. In particular, we think Metzinger is right to stress
215 both the value of first-person experiences and to eschew the idea that the *content* of these
216 experiences is of primary explanatory value.

217

218 Instead, we think that the primary value of first-person experience is best considered in terms
219 of its effect on the meta-problem. Sometimes first-person demonstrations *that* something can
220 be done are more powerful than *what* is actually accomplished: the first-person value of
221 direct interventions, we suspect, will be most useful in that regard. Indeed, that usefulness
222 might obtain even before we have a full science of consciousness sorted—for although the
223 meta-problem and the hard problem are intertwined, it may take less to fix the former than it
224 will to solve the latter.

225

226 Ultimately, we think Chalmers’ framework should be of great interest to the experimentalist
227 and the philosopher alike. So long as the meta-problem remains pressing, experimental

228 research will seem unsatisfying. Conversely, experimental research might itself hold the key
229 to fixing the meta-problem, and thereby making real progress on the science of
230 consciousness.

231

232 *Acknowledgements*

233 Thanks to Tim Bayne, Jakob Hohwy, and Bjorn Merker for discussion, and to François
234 Kammerer and an anonymous referee for helpful comments on a previous draft. This paper
235 was funded by Australian Research Council Grant FT140100422 (to CK) and FT140100452
236 (to ABB).

237

238 **References:**

- 239 Armstrong D (1999) *The Mind-body Problem: An Opinionated Introduction*. Westview
240 Press: Boulder.
- 241 Barron, A.B. and Klein, C. (2016) What insects can tell us about the origins of consciousness.
242 *Proceedings of the National Academy of Sciences* 113(18): 4900-4908.
- 243 Campbell J (2007) An interventionist approach to causation in psychology. in *Causal*
244 *Learning: Psychology, Philosophy and Computation*, ed Alison Gopnik and Laura
245 Schulz. Oxford: Oxford University Press, pp 58-66.
- 246 Chalmers D (1996) *The Conscious Mind: In Search of a Fundamental Theory* (Oxford
247 University Press, New York).
- 248 Chalmers, D. (2018) The Meta-Problem of Consciousness. *Journal of Consciousness Studies*
249 25(9-10): 6-61.
- 250 Feigl, H. (1958) "The 'mental' and the 'physical'" *Minnesota studies in the philosophy of*
251 *science*. 2(2): 370-497.
- 252 Hilbert DR (1987) *Color and color perception: A study in anthropocentric realism* (Center
253 for the Study of Language and Information, Stanford).
- 254 Hitchcock CR (1996) The Role Of Contrast In Causal And Explanatory Claims. *Synthese*
255 107:395-419.
- 256 Jackson F (1982) Epiphenomenal qualia. *The Philosophical Quarterly* 32:127-136.
- 257 Klein, C. and Barron, A.B. (ms in review) How Experimental Neuroscientists Can Fix the
258 Hard Problem of Consciousness.
- 259 Klein, C. (2017) Brain Regions as Difference-Makers. *Philosophical Psychology* 30(1-2): 1-
260 20.
- 261 Langlitz (2016) Is There A Place For Psychedelics In Philosophy? *Common Knowledge* 22:3:
262 373-384.
- 263 Lashley KS (1960) *The neuropsychology of Lashley: Selected papers of KS Lashley*. ed F. A.
264 Beach, D. O. Hebb, C. T. Morgan & H. W. Nissen. New York: McGraw-Hill.
- 265 Metzinger T (2006) "Reply to Hobson: Can There Be a First-Person Science of
266 Consciousness?" *Psyche* 12, no. 4: 2.
- 267 Pettit P (2003) Looks as powers. *Philosophical Issues* 13(1):221-252.
- 268 Ramberg (2000) The death of vitalism and the birth of organic chemistry: Wohler's urea
269 synthesis and the disciplinary identity of organic chemistry. *Ambix* 47(3): 170-195.
- 270 Salmon, W. (1989) *Four Decades of scientific explanation*. Minneapolis: University of
271 Minneapolis press.
- 272 Woodward J (2003) *Making Things Happen* (Oxford University Press, New York).

273 Woodward, J (2010) Causation in biology: Stability, specificity, and the choice of levels of
274 explanation. *Biology and Philosophy* 25(3): 287-318.
275
276