

1 EVERYTHING AND MORE: THE PROSPECTS OF WHOLE BRAIN 1  
 2 EMULATION\* 2

3  
 4 **W**hole brain emulation (WBE) has been proposed as the 3  
 4 most promising avenue for creating human-level artificial 4  
 5 intelligence, creating superintelligence, and even for 5  
 6 achieving immortality.<sup>1</sup> The basic goal behind WBE is to create a soft- 6  
 7 ware model of one's mind which could then be uploaded to a stor- 7  
 8 age system. The resulting software model could then be downloaded 8  
 9 to new hardware, and the uploaded model could, so the idea goes, 9  
 10 be used to recreate a functional isomorph of the original brain from 10  
 11 which it was copied. Doing so would then, by hypothesis, replicate all 11  
 12 the psychological features that were present in the original individual 12  
 13 whose brain was copied. 13

14 What level of detail might be needed for an upload? *Prima facie*, 14  
 15 all one would need is a "connectome."<sup>2</sup> A connectome is an anatomical 15  
 16 wiring diagram that charts the connections between each neuron, 16  
 17 giving us an overall wiring map of the brain. The idea is that by repli- 17  
 18 cating a wiring diagram of one's brain, we would thereby replicate 18  
 19 one's psychology.<sup>3</sup> 19

20 Connectomics and WBE are natural partners, and together they ap- 20  
 21 pear to offer tantalizing possibilities. Connectomics is a well-defined 21  
 22 research program, one which is well underway: The Human Con- 22  
 23 nectome Project is an interinstitutional, reputable, amply funded re- 23  
 24 search project.<sup>4</sup> The success of the project could itself underwrite 24  
 25 WBE, as the wiring diagram of the connectome may be thought to 25  
 26 replicate the functional characteristics of one's brain. Consequently, 26  
 27 WBE seems particularly well placed among all transformative tech- 27  
 28 nologies as, crucially, it need not rely on any conceptual breakthrough 28  
 29 29  
 30 30

31 \* Helpful comments and discussion were received from Ron Avni, Adam Bradley, 31  
 32 David Chalmers, Tim Crane, Cian Dorr, Jackson Kernion, Jessica Moss, Jake Quilty- 32  
 33 Dunn, Shen Pan, Kate Ritchie, David Rosenthal, Fiona Schick, Henry Schiller, and 33  
 34 David Udell. They are all hereby thanked for their collegiality, patience, and friend- 34  
 35 ship. 35

36 <sup>1</sup> Anders Sandberg and Nick Bostrom, *Whole Brain Emulation: A Roadmap*, Technical 36  
 Report 2008-3, Future of Humanity Institute, Oxford University, 2008, [www.fhi.ox.ac.uk/reports/20083.pdf](http://www.fhi.ox.ac.uk/reports/20083.pdf). 36

37 <sup>2</sup> O. Sporns, G. Tononi, and R. Kötter, "The Human Connectome: A Structural De- 37  
 38 scription of the Human Brain," *PLoS Computational Biology*, 1, 4 (2005): e42. 38

39 <sup>3</sup> Sebastian Seung, *Connectome: How the Brain's Wiring Makes Us Who We Are* (New York: 39  
 Houghton Mifflin Harcourt, 2012). 39

40 <sup>4</sup> See <http://www.humanconnectomeproject.org/> for details. 40  
 41 41

1 to ensure its success.<sup>5</sup> Here are Anders Sandberg and Nick Bostrom, 1  
2 two prominent futurists, on WBE's promise: 2

3 WBE represents a formidable engineering and research problem, yet 3  
4 one which appears to have a well-defined goal and could, it would seem, 4  
5 be achieved by extrapolations of current technology. This is unlike many 5  
6 other suggested radically transformative technologies like artificial intel- 6  
7 ligence where we do not have any clear metric of how far we are from 7  
8 success.<sup>6</sup> 8

9 Sandberg and Bostrom suggest that WBE allows for a "clear metric" 9  
10 because we can understand how far we are from replicating a complete 10  
11 brain's wiring diagrams. In this way WBE stands alone as the only 11  
12 route to superintelligence that we currently appear to understand, at 12  
13 least in broad strokes. Moreover, we have made progress on this front, 13  
14 with some simple species' connectomes already mapped (for exam- 14  
15 ple, *C. Elegans*'s connectome was mapped over thirty years ago).<sup>7</sup> 15

16 At first blush, WBE seems tantalizing. Understanding the brain and 16  
17 mind is far too difficult a task to accomplish in any reasonable amount 17  
18 of time. WBE however, holds promise of being able to sidestep this 18  
19 worry: "A key assumption, characteristic of the WBE approach to AI, 19  
20 is nonorganicism: total understanding of the brain is not needed, 20  
21 just understanding of the component parts and their functional in- 21  
22 teractions."<sup>8</sup> WBE holds out hope that we can emulate the brain's 22  
23 functional apparatus without understanding, for example, how neu- 23  
24 ral structure itself leads to intentionality, consciousness, or intelli- 24  
25 gence. Because of how successful connectomics has been, Sandberg 25  
26 estimates a 50% confidence level in the proposition that WBE will 26  
27 arise by 2064.<sup>9</sup> Even the rosiest optimists among us would need to put 27  
28 the chances that we have achieved anywhere near full understanding 28  
29 of the brain, never mind the mind, by then as infinitesimal. Thus, 29  
30 WBE seems much more promising than, say, creating complete mod- 30  
31 els of the mind based on understanding everything about our dis- 31  
32 parate mental faculties and capacities (for example, our characters, 32  
33 personalities, or ways of acquiring beliefs). 33

34  
35  
36 <sup>5</sup>This claim assumes that we would not need conceptual breakthroughs in neuro- 35  
36 science to understand the connectome. 36

37 <sup>6</sup>Sandberg and Bostrom, *Whole Brain Emulation*, *op. cit.*, p. 8. 37

38 <sup>7</sup>J. White et al., "The Structure of the Nervous System of the Nematode *Caenorhab- 38  
39 ditis Elegans*," *Philosophical Transactions of the Royal Society London B: Biological Sciences*, 39  
ccciv, 1165 (1986): 1–340. 39

40 <sup>8</sup>Anders Sandberg, "Feasibility of Whole Brain Emulation," in Vincent C. Müller, 40  
41 ed., *Philosophy and Theory of Artificial Intelligence* (Berlin: Springer, 2013), pp. 251–64, at 41  
42 p. 257. 42

43 <sup>9</sup>Anders Sandberg, "Monte Carlo Model of Brain Emulation Development," Work- 43  
ing Paper 2014-1, Future of Humanity Institute, Oxford University, 2014. 43

1 WBE is also tantalizing because it allows for the possibility of ex- 1  
 2 tremely lofty goals. If reproducing the connectome would reproduce 2  
 3 the functional properties of the mapped brain, then WBE might hold 3  
 4 the key for immortality. One's identity might be thought of as the 4  
 5 totality of one's psychology—one's personality, memories, emotions, 5  
 6 and the like. But the promise of WBE allows for other lofty goals be- 6  
 7 sides immortality, as WBE also seems like the clearest route to achiev- 7  
 8 ing superintelligence.<sup>10</sup> WBE would allow for relatively cheap upload- 8  
 9 ing and storage of human-level intelligence (which itself would consti- 9  
 10 tute "weak superintelligence," human-level intelligence that can oper- 10  
 11 ate at much greater speeds).<sup>11</sup> As human capital is the central driver of 11  
 12 economic growth, having large amounts of readily available human- 12  
 13 level intelligences will make for enormous technological and societal 13  
 14 enhancement.<sup>12</sup> 14

15 Since few technologies hold the promise of such transformative 15  
 16 ends as immortality and superintelligence, the question of the feasibil- 16  
 17 ity of WBE is pressing, even if one's *a priori* intuitions of the chances of 17  
 18 achieving it are more pessimistic than Sandberg's. My goal in this pa- 18  
 19 per is to provide that analysis. I argue that one should have a healthy 19  
 20 skepticism as to the fecundity of WBE. Moreover, the problems with 20  
 21 WBE are not specific to it—showing the problems inherent in WBE 21  
 22 will illuminate fissures in the doctrine of computationalism writ large. 22

### 23 I. WHOLE MIND EMULATION 23

24 I am interested in two questions. First, would a connectome of a 24  
 25 single mind suffice to instantiate a broad range of psychological fea- 25  
 26 tures, features such as one's personality, character, intelligence, and 26  
 27 phenomenology?<sup>13</sup> Second, assuming an affirmative response to the 27  
 28 first question, would the connectome suffice for personal identity? 28  
 29 29  
 30 30  
 31 31

32 <sup>10</sup>Superintelligence is generally glossed as intelligence that far exceeds human ca- 32  
 33 pacities in every domain; see Nick Bostrom, *Anthropic Bias: Observation Selection Effects*  
 34 *in Science and Philosophy* (New York: Routledge, 2010). It is unclear how much a system  
 35 must exceed human intelligence in order to qualify as superintelligent. 35

36 <sup>11</sup>Nick Bostrom, "How Long before Superintelligence?," *Linguistic and Philosophical*  
 37 *Investigations*, v, 1 (2006): 11–30. 36

37 <sup>12</sup>See Robin Hanson, "Economics of the Singularity," *IEEE Spectrum*, (Jun. 1, 2008):  
 38 37–42; Robin Hanson, "If Uploads Come First: The Crack of a Future Dawn," *Entropy*,  
 39 vi, 2 (1994): 10–15; and Robin Hanson, *The Age of Em: Work, Love, and Life when Robots*  
 40 *Rule the Earth* (Oxford: Oxford University Press, 2016). 37

38 <sup>13</sup>A connectome is a representation, so it may be better to speak of an instantiation  
 39 of a connectome. For convenience I will speak as if connectomes are instantiations, as  
 40 the distinction will not affect my argument. Similarly, when I speak of uploading con-  
 41 nectomes, one may prefer to think about uploading instantiations of connectomes. See  
 42 Susan Schneider, *Artificial You: AI and the Future of Your Mind* (Princeton, NJ: Princeton  
 43 University Press, 2019), chapter 8. 39  
 40  
 41  
 42  
 43

1 Would the “software” to be uploaded ensure duplicating a given 1  
 2 mind? Would uploading my connectome suffice for uploading *me*?<sup>14</sup> 2

3 To put it in Sandberg and Bostrom’s terms, assume a brain emu- 3  
 4 lator is a piece of software. Our question is then whether brain emu- 4  
 5 lation so understood would entail *mind emulation*—a model that is 5  
 6 “detailed and correct enough to produce the phenomenological ef- 6  
 7 fects of a mind” *inter alia*.<sup>15</sup> 7

8 WBE can be a success even if it cannot preserve full personal iden- 8  
 9 tity. What it means to be successful depends on one’s aims: if a connec- 9  
 10 tome suffices for establishing propositional attitudes, then uploading 10  
 11 a connectome could allow for creating artificial intelligence.<sup>16</sup> If this 11  
 12 intelligence can be harnessed, then WBE might serve as the catalyst 12  
 13 for superintelligence, even if WBE could not ensure immortality. 13

## 14 II. PHYSICALISM, MULTIPLE REALIZABILITY, AND WBE 14

15 WBE seems *prima facie* feasible. It is a natural bedfellow of computa- 15  
 16 tionalism, the idea that the mind is just a computer of sorts, where 16  
 17 mental processes are understood as transformations of mental repre- 17  
 18 sentations.<sup>17</sup> In a canonical formulation of computationalism, mental 18  
 19 representations are taken to be symbols, and mental processes com- 19  
 20 pute over the formal properties of the symbols.<sup>18</sup> Computationalism 20  
 21 gains inspiration from the Church-Turing thesis, which holds that any 21  
 22 computable function can be computed by a Turing machine. Any soft- 22  
 23 ware duplicate will be Turing-machine equivalent. Unless one believes 23  
 24 that the mind is somehow outside of the physical realm altogether, 24  
 25 there should be no *a priori* restriction to the feasibility of WBE. 25  
 26 26

27  
 28  
 29 <sup>14</sup>Or would uploading my connectome suffice for even uploading a token of the 29  
 30 type that is me? See *ibid.* for discussion. 30

31 <sup>15</sup>Sandberg and Bostrom, *Whole Brain Emulation*, *op. cit.*, p. 7. 31

32 <sup>16</sup>Strictly speaking, propositional attitudes may not be necessary for intelligence. Per- 32  
 33 haps an agent could be intelligent without believing or desiring anything. And strictly 33  
 34 speaking, attitudes may not be sufficient either: maybe there could be a creature with 34  
 35 beliefs and desires but no combinatorial apparatus for generating rational thought or 35  
 36 behavior. That said, having full-blooded propositional attitudes seems to make intelli- 36

37 <sup>17</sup>Schneider, *Artificial You*, *op. cit.*, argues against the canonical reading of computa- 37  
 38 tionalism while still holding a broadly computationalist theory. 38

39 <sup>18</sup>See, for example, Jerry A. Fodor, *Concepts: Where Cognitive Science Went Wrong* (Ox- 39  
 40 ford: Oxford University Press, 1998); and Ned Block, “The Mind as the Software of the 40  
 41 Brain,” in D. Osherson and E. Smith, eds., *Thinking: An Invitation to Cognitive Science*, 41  
 42 *Volume 3*, 2nd ed. (Cambridge, MA: MIT Press, 1995), pp. 377–426. Of course, some 42  
 43 other theorists that might be considered broadly computationalist—for example, cer- 43  
 44 tain connectionists—would balk at discussions of symbolic computation. For a skeptical 44  
 45 take on classical computationalist models, see Schneider, *Artificial You*, *op. cit.*; and Susan 45  
 46 Schneider, *The Language of Thought: A New Philosophical Direction* (Cambridge, MA: 46  
 47 MIT Press, 2011). 47

1 To make the case as strong as possible for WBE, let us assume token 1  
2 physicalism in what follows so that the Church-Turing thesis holds 2  
3 over all neural events.<sup>19</sup> Given that, what are the roadblocks to WBE? 3

4 Sandberg and Bostrom write that WBE only requires that we find 4  
5 “a 1-to-1 model where all relevant properties of a system exist.”<sup>20</sup> But 5  
6 what are the relevant properties of one’s brain? This question is press- 6  
7 ing. WBE relies on the idea that the end goal of computational neuro- 7  
8 science is to provide a neuroinformatic map of the brain. The de- 8  
9 tail of the map matters: if WBE requires a level of detail equivalent 9  
10 to molecule-for-molecule duplication, then it is far too information 10  
11 rich a plan to be feasible in the short term (where the short term 11  
12 includes times as early as Sandberg’s estimate of 2064). WBE propo- 12  
13 nents understand that they need to abide by “nonorganicism,” and 13  
14 instead suggest that only a certain level of functional understanding 14  
15 should be necessary for WBE to be viable. For example, Sandberg 15  
16 writes, “For the current paper we will focus on simulations that at- 16  
17 tempt to achieve full functional equivalence—all relevant behavioral 17  
18 properties and internal causal links of the original system are repli- 18  
19 cated.”<sup>21</sup> Part of this project entails modeling the interactions of “neu- 19  
20 rons and brain systems, and the emergent dynamics between them.”<sup>22</sup> 20  
21 But what are the *brain systems* referred to here? If they are merely 21  
22 wiring diagrams between neurons, then although mapping out the 22  
23 connections is an extremely difficult engineering task, it nonetheless 23  
24 is one that seems feasible. However, if more than the connectome 24  
25 matters, if, instead, lower-level, finer-grained details, such as ones that 25  
26 involve neurochemical elements, or other substances that correspond 26  
27 to our “hardware,” are germane, then the road to emulation is much 27  
28 less clear. The problem in front of us is to identify whether there are 28  
29 relevant aspects of cognition broadly construed (phenomenology, in- 29  
30 tentionality, intelligence, and personality, at a minimum) that are not 30  
31 merely dictated by the connectome. 31

32 Would the connectome suffice for replicating functional compe- 32  
33 tence? To answer affirmatively is to presuppose a version of machine 33  
34 34  
35 35

36 <sup>19</sup>Token physicalism itself might be too strong a thesis for some—many non-dualists 36  
37 and most property dualists seem to reject it. The *a priori* feasibility of WBE only needs 37  
38 a weaker thesis, something such as the mental supervening on the physical. Never- 38  
39 theless, assuming token physicalism will not matter much beyond helping explication 39  
40 (and will make the WBE proponent’s case stronger). For more on token physicalism, 40  
41 see Daniel Stoljar, “Physicalism,” in Edward N. Zalta, ed., *The Stanford Encyclopedia of* 41  
*Philosophy* (Winter 2021 Edition), <https://plato.stanford.edu/entries/physicalism/>. 41

42 <sup>20</sup>Sandberg and Bostrom, *Whole Brain Emulation*, *op. cit.*, p. 7. 42

43 <sup>21</sup>Sandberg, “Feasibility of Whole Brain Emulation,” *op. cit.*, p. 253. 43

<sup>22</sup>*Ibid.*, p. 252. 43

1 functionalism<sup>23</sup> as well as the multiple realizability thesis,<sup>24</sup> the idea 1  
 2 that psychological properties can be realized from a wide array of 2  
 3 structural properties.<sup>25</sup> Both views are closely related: machine func- 3  
 4 tionalism dictates that all essential properties of the mind are func- 4  
 5 tional (and not structural) properties;<sup>26</sup> that is, it assumes that, for 5  
 6 example, to be a belief is to just be a mental state that serves a certain 6  
 7 role, the role that belief generally serves.<sup>27</sup> This function is the essence 7  
 8 of the mental state. As such, there is nothing in the essence of belief, 8  
 9 for example, that makes it seem as if it had to be realized by a partic- 9  
 10 ular substrate. Thus, perhaps there could be intelligent creatures that 10  
 11 had silicon “brains.” If these creatures had beliefs, this would prove 11  
 12 that belief is multiply realizable, since it could be realized in brains 12  
 13 like ours or in heads filled with silicon. 13  
 14

### 15 III. PROBLEMS FOR WHOLE BRAIN EMULATION 15

16 Let us start with a seemingly pressing, though relatively easy, problem 16  
 17 for the multiple realizability thesis: the embodied mind and extended 17  
 18 mind theses. These theories hold that our minds extend beyond our 18  
 19 skulls, and not just because of, for example, an externalist semantics 19  
 20 that dictates that content is not only in the head. Proponents of the 20  
 21 embodied mind posit that the body is integral to the functioning of 21  
 22 the mind. Locomotion is interpreted as a central cognitive function, 22  
 23 not one that is just useful for aiding in cognitive development but 23  
 24 instead partially constitutive of cognition itself. Similarly, extended 24  
 25 25  
 26  
 27  
 28

29 <sup>23</sup> Hilary Putnam, “The Nature of Mental States,” in *Philosophical Papers, Volume 2: Mind, Language, and Reality* (Cambridge, UK: Cambridge University Press, 1975). 29

30 <sup>24</sup> J. A. Fodor, “Special Sciences (Or: The Disunity of Science as a Working Hypothe- 30  
 31 sis),” *Synthese*, xxviii, 2 (1974): 97–115. 31

32 <sup>25</sup> This is close to right, though there is a bit of slippage. The former thesis—whether 32  
 33 the connectome would suffice for replicating functional competence—is about behav- 33  
 34 iorally competence, whereas the latter—machine functionalism—is about the essence of 34  
 35 the mind. All of those who answer the latter question affirmatively will do the same for 35  
 36 the former, but some who answer the former affirmatively may be silent on the latter. 36

37 <sup>26</sup> Which properties count as structural—say, the fusiform gyrus, or an electron— 37  
 38 depends on one’s explanatory ends. Properties that look structural from one vantage 38  
 39 point (for example, the prefrontal cortex from the standpoint of intentional psychol- 39  
 40 ogy) look functional from another (for example, the prefrontal cortex from the stand- 40  
 41 point of biochemistry). See William G. Lycan, *Consciousness* (Cambridge, MA: MIT 41  
 42 Press, 1987). 42

43 <sup>27</sup> Ironically, some of the biggest proponents of connectomics also hold that neural 43  
 44 structure and function are closely related, resulting in a rather precarious dialectic 44  
 45 position (see, for example, Beth L. Chen, David H. Hall, and Dmitri B. Chklovskii, 45  
 46 “Wiring Optimization Can Relate Neuronal Structure and Function,” *Proceedings of the 46  
 47 National Academy of Sciences*, ciii, 12 (2006): 4723–28). 47

1 mind theorists claim that objects outside of one's body entirely—say 1  
2 one's cellphone—partially constitute one's cognitive apparatus.<sup>28</sup> 2

3 Both the embodied and extended mind theses seem in tension with 3  
4 the multiple realizability thesis, which presupposes that one can up- 4  
5 load the cognitive software into any number of hardware realizers. But 5  
6 if embodied and extended cognition theorists are right, then there 6  
7 are real restrictions on the types of programs one could be uploaded 7  
8 into—for instance, a brain in a vat would not suffice for cognition. 8

9 Nonetheless, the embodied mind is not a deep obstacle for the fea- 9  
10 sibility of WBE. For one thing, the restrictions that would apply are, 10  
11 in the scheme of things, relatively trivial—they are not restrictions on 11  
12 the type of hardware that would be needed for uploads, but are in- 12  
13 stead restrictions on the type of environments in which the hardware 13  
14 would have to be embedded. Adding the analogs of perceptual inputs 14  
15 and motor outputs, as well as some objects to interact with, is far less 15  
16 challenging than successfully reproducing an entire functional copy 16  
17 of a brain. In the case where we are envisioning that we can already 17  
18 do the latter, the former should be a small roadblock at most. 18

19 However, there is a more serious problem lurking, one that ques- 19  
20 tions the scope of functionalism in its entirety. Functionalism about 20  
21 mental states—beliefs, desires, hope, and the like—seems appealing 21  
22 because the functional role that each state plays seems essential to its 22  
23 character. For example, if you found a state that was not caused by 23  
24 perception, did not interact with motivational states to produce be- 24  
25 havior, and did not serve as premises in inferences then it would be 25  
26 hard to see how it could count as a belief.<sup>29</sup> Though that functional- 26  
27 ist intuition is reasonable enough, extending it to other states—say, 27  
28 phenomenological and motivational states—is a much more tenuous 28  
29 proposition. Beliefs seem functional, but, for example, experiencing 29  
30 something as green seems less so.<sup>30</sup> 30

31 Thus we again face the question: what level of granularity is nec- 31  
32 essary in order for WBE? What properties are the relevant ones that 32  
33 need to be recreated in order for emulation to be successful? One's 33  
34 take on whether the connectome will recapitulate cognition writ large 34  
35 will depend on one's theory of consciousness, even among physical- 35  
36 ists. 36  
37 37  
38 38

39 <sup>28</sup> Andy Clark and David Chalmers, "The Extended Mind," *Analysis*, LVIII, 1 (1998): 39  
40 7–19. 40

41 <sup>29</sup> Jake Quilty-Dunn and Eric Mandelbaum, "Against Dispositionalism: Belief in Cog- 41  
42 nitive Science," *Philosophical Studies*, CLXXV, 9 (2018): 2353–72. 42

43 <sup>30</sup> See also Tim Maudlin, "Computation and Consciousness," this JOURNAL, LXXXVI, 43  
44 8 (1989): 407–32. 44

1 The theories of consciousness that matter for evaluating the pros- 1  
2 pects of WBE are theories of what makes a state phenomenally con- 2  
3 scious. There are three major ones: (1) the higher-order thought 3  
4 (HOT) theory; (2) the global workspace theory; and (3) the biologi- 4  
5 cal theory. 5

6 All these theories give different answers to the question of what 6  
7 makes a state conscious. The HOT theory states that a first-order men- 7  
8 tal state becomes phenomenally conscious when a higher-order state 8  
9 takes the first state as its content.<sup>31</sup> In this way, consciousness is re- 9  
10 duced to thoughts about mental states: a thought becomes conscious 10  
11 when we have another thought about it; a feeling becomes conscious 11  
12 when we have a thought about that feeling; and so on. HOT the- 12  
13 ory is friendly to WBE because the essential relation posited—that 13  
14 of a thought monitoring another mental state—is a functional one.<sup>32</sup> 14  
15 Monitoring does not involve implementation or machinery at all, so 15  
16 it should be able to be instantiated in many different ways, and thus 16  
17 ought to be amenable to multiple realizability. 17

18 Global workspace theory posits that any state that is “globally broad- 18  
19 cast” is *ipso facto* phenomenally conscious. To be globally broadcast is 19  
20 a dispositional property: it is to be a state that is ready to be utilized by 20  
21 a varied array of other mental processes, such as reasoning, linguistic, 21  
22 and motoric processes.<sup>33</sup> Popular versions of global broadcasting posit 22  
23 competitive neural networks where sensory and frontal areas compete 23  
24 for resources, with the winner becoming conscious.<sup>34</sup> But none of the 24  
25 neural details, not even the use of neural nets, are essential to the 25  
26 view; instead, they are just one way to flesh it out. At its core, the 26  
27 global workspace theory is a functionalist view: it hypothesizes that 27  
28 to be phenomenally conscious is just to be a representation that is 28  
29 available to a host of consuming mental mechanisms (for example, 29  
30 language production). In us, such consumption may involve details 30  
31 about competition between sensory and frontal cortices; however, that 31  
32 is a detail (most likely) about consciousness in us and not phenome- 32  
33 nal consciousness *simpliciter*. Thus, global workspace theory allows that 33  
34 if we uploaded our connectome, we could replicate consciousness in a 34  
35 35  
36 36

37 <sup>31</sup>David Rosenthal, *Consciousness and Mind* (New York: Oxford University Press, 37  
38 2005). 38

39 <sup>32</sup>One could, in theory, be a HOT theorist but not be a functionalist (if, say, one 39  
40 rejected functionalism about intentionality and had a non-functional specification of 40  
41 monitoring). That said, I cannot think of any non-functionalist HOT theorists. 40

41 <sup>33</sup>Bernard J. Baars, *A Cognitive Theory of Consciousness* (Cambridge, UK: Cambridge 41  
42 University Press, 1988). 42

43 <sup>34</sup>Stanislas Dehaene et al., “Conscious, Preconscious, and Subliminal Processing: A 43  
44 Testable Taxonomy,” *Trends in Cognitive Science*, x, 5 (2006): 204–11. 44



1 nonbiological substrate. In other words, global workspace, like HOT, 1  
2 is an essentially functionalist theory, one compatible with the multiple 2  
3 realizability thesis and WBE. 3

4 However, the third theory of consciousness—the biological theory, 4  
5 —is where deep problems for WBE arise. The biological theory 5  
6 posits that the coding and interchange of information between elec- 6  
7 trical and chemical formats gives rise to consciousness, and that the 7  
8 specific neural hardware we use is essential to phenomenal conscious- 8  
9 ness.<sup>35</sup> 9

10 Some prominent arguments for the biological theory come from 10  
11 Ned Block: one argument relies on the explanatory gap, and the 11  
12 other relies on perceptual overflow. The explanatory gap is the 12  
13 thesis that we have no idea how a subjective state (such as seeing red or 13  
14 hearing middle C on a piano) could be identical to an objective state 14  
15 (such as having a certain pattern of neuronal activation).<sup>36</sup> The thesis 15  
16 does not claim that humans cannot in principle explain how objective 16  
17 states could give rise to subjective states. Instead, it is a theory about 17  
18 our current epistemic position, one which claims that at this moment 18  
19 we have no clue how psychophysical identities could be true.<sup>37</sup> The 19  
20 idea is that we do not yet possess the concepts to bridge this gap (al- 20  
21 though one day we may). 21  
22

23 While few dispute the explanatory gap's existence, the morals one 23  
24 should draw from it are more controversial. The biological theory 24  
25 takes the existence of the explanatory gap as support, as neither the 25  
26 HOT nor the global workspace theory can explain why, if conscious- 26  
27 ness is a functional property, we should have an explanatory gap and 27  
28 the subsequent “hard problem.”<sup>38</sup> 28  
29

30 Another argument Block puts forward in favor of the biological theory 30  
31 is that it is the only view that can explain phenomenal overflow. 31  
32 “Phenomenal overflow” describes situations where one's phenomenal 32  
33 consciousness—generally in perceptual situations— 33  
34

35 Ned Block, “Comparing the Major Theories of Consciousness,” in M. Gazzaniga, 37  
38 ed., *The Cognitive Neurosciences* (Cambridge, MA: MIT Press, 2009), pp. 1111–22. 38

36 Joseph Levine, “Materialism and Qualia: The Explanatory Gap,” *Pacific Philosophi- 39  
cal Quarterly*, LXIV, 4 (1983): 354–61. 39

37 *Ibid.* 40

38 Block, “Comparing the Major Theories of Consciousness,” *op. cit.* I do not quite 41  
42 see how the explanatory gap is supposed to help the biological theory here, as it seems 42  
43 to also fall prey to the gap. The argument is presented here out of completeness, not 43  
44 endorsement. 44

1 overflows cognitive access. Perception and phenomenal conscious- 1  
 2 ness more generally seem *richer* than what cognition can conceptual- 2  
 3 ize. Parade examples use the logic of a partial report paradigm.<sup>39</sup> Sub- 3  
 4 jects see an arrangement of letters (for example, three rows of four 4  
 5 letters) for a brief period of time. The letters then disappear and the 5  
 6 subjects are cued to one of the rows. Subjects can report three or four 6  
 7 letters from any cued row. But if subjects are asked to report as many 7  
 8 letters as possible without any cue, they can still report only three or 8  
 9 four letters. That is, subjects appear to consciously see all of the let- 9  
 10 ters during the presentation but can only consciously access three or 10  
 11 four total letters from the twelve-letter array. The rest of the letters are 11  
 12 consciously perceived—they add to one’s phenomenology—but they 12  
 13 are not consciously accessed.<sup>40</sup> 13

14 Theorists like Block use overflow to argue for the biological the- 14  
 15 ory.<sup>41</sup> They argue that any functional view of consciousness, such as 15  
 16 HOT or global workspace, would place the unseen letters in the app- 16  
 17 ropriate functional role as dictated by those theories. Take global 17  
 18 workspace theory: the letters are originally conscious (because they 18  
 19 add to phenomenology); since they are conscious, they should be re- 19  
 20 reportable because, by hypothesis, to be conscious is just to be available 20  
 21 in the workspace, which entails being available to report. But the let- 21  
 22 ters are not reportable even though they are conscious; thus, Block 22  
 23 reasons, the global workspace theory must be wrong. 23

24 The biological theory of consciousness is the only non-functional- 24  
 25 istic of the theories canvassed, and as such it can explain the richness 25  
 26 of perception and experience by interpreting that richness as overflow- 26  
 27 ing access. What makes a state conscious is not its dispositional prop- 27  
 28 erties (for example, being available to report or being the content of 28  
 29 another thought) but merely the state being caused (or realized) by 29  
 30 the specific biological machinery we have. 30

31 The biological theory also finds support outside of any of the over- 31  
 32 flow arguments. The connectome is the level of grain that most the- 32  
 33 orists find plausible for positing as the functional basis of the mind.<sup>42</sup> 33  
 34 34  
 35 35

36 <sup>39</sup> George Sperling, “The Information Available in Brief Visual Presentations,” *Psy-* 36  
 37 *chological Monographs*, LXXIV, 11 (1960): 1–29. 37

38 <sup>40</sup> For competing takes on overflow, see Ian Phillips, “No Watershed for Overflow: 38  
 39 Recent Work on the Richness of Consciousness,” *Philosophical Psychology*, xxix, 2 (2016): 39  
 40 236–49; and Steven Gross and Jonathan Flombaum, “Does Perceptual Consciousness 40  
 41 Overflow Cognitive Access? The Challenge from Probabilistic, Hierarchical Processes,” 41  
 42 *Mind and Language*, xxxii, 3 (2017): 358–91. 42

43 <sup>41</sup> Ned Block, “Perceptual Consciousness Overflows Cognitive Access,” *Trends in Cog-* 43  
 44 *nitive Science*, xv, 12 (2011): 567–75. 44

<sup>42</sup> Seung, *Connectome*, *op. cit.*

1 But the connectome is just an anatomical wiring diagram—even elec- 1  
 2 trical connections between neurons are left out.<sup>43</sup> *A fortiori*, connec- 2  
 3 tomics is committed to the view that a sub-neuronal difference should 3  
 4 not lead to a functional difference. But sub-neuronal differences do 4  
 5 appear to lead to psychological differences. What causes the vast in- 5  
 6 dividual differences in phenomenology is extremely unclear at the 6  
 7 moment. But the contribution of sub-neuronal properties is inte- 7  
 8 gral in a way that is rarely appreciated in the literature. Serotonin, 8  
 9 dopamine, norepinephrine, histamine, and countless neuropeptides 9  
 10 are not accounted for in the connectome; they count as part of the 10  
 11 “hardware” of our system. These neurochemical properties act as neu- 11  
 12 romodulators, affecting neuronal connections in fundamental ways, 12  
 13 even changing basic neuronal functions.<sup>44</sup> In the connectome of *C* 13  
 14 *Elegans*—a vastly easier connectome to understand than the human 14  
 15 one—every neuron and synapse was subject to neuromodulation.<sup>45</sup> 15  
 16 The effect of neuromodulation is enormous in all nervous systems:<sup>46</sup> 16  
 17

18 Modulators can qualitatively alter the neuron’s intrinsic properties, 18  
 19 transforming neurons from tonic spiking to those generating plateau 19  
 20 potentials or bursts. The effect of neuromodulators can activate or si- 20  
 21 lence an entire circuit, change its frequency, and/or the phase relation- 21  
 22 ships of the motor patterns generated.<sup>47</sup> 22

23 And again, this holds in creatures much simpler than human beings 23  
 24 (for example, in worms); it is reasonable to suppose that in the more 24  
 25 baroque case of the human brain, neuromodulators (to say noth- 25  
 26 ing of glial cells) take on an even greater role. After all, we depend 26  
 27 on intervening on neuromodulators to change affective and motiva- 27  
 28 tional states—serotonin reuptake pharmaceuticals are not targeting 28  
 29 neuronal connections but neurochemicals. To put it mildly, it seems 29  
 30 implausible that every neuron can have its basic function changed by 30  
 31 its instantiation base yet also hold that the instantiation base would 31  
 32 have no effect on any cognitive property. 32  
 33  
 34  
 35

36 <sup>43</sup>Joshua L. Morgan and Jeff W. Lichtman, “Why Not Connectomics?,” *Nature Meth-*  
 37 *ods*, x, 6 (2013): 494–500.

38 <sup>44</sup>Cornelia I. Bargmann and Eve Marder, “From the Connectome to Brain Func-  
 39 *tion*,” *Nature Methods*, x, 6 (2013): 483–90.

40 <sup>45</sup>*Ibid.* *C. Elegans*’s connectome is perhaps not the most favorable piece of evidence  
 41 for WBE enthusiasts. It was mapped in 1986, and yet we still have little idea what func-  
 42 tion any of its neuronal connections subserves, even though it only has about 300 neu-  
 43 rons as opposed to our 100 billion or so neurons.

44 <sup>46</sup>Eve Marder, “Neuromodulation of Neuronal Circuits: Back to the Future,” *Neuron*,  
 45 LXXVI, 1 (2012): 1–11.

46 <sup>47</sup>Bargmann and Marder, “From the Connectome to Brain Function,” *op. cit.*, p. 486.

1 Moreover, we have good evidence that some sub-connectomic prop- 1  
 2 erties do matter for psychology. For instance, steroids from the 2  
 3 adrenal cortex, as well as from sex organs, are not captured by the 3  
 4 connectome.<sup>48</sup> But increases in (for example) testosterone plainly do 4  
 5 affect a wide range of behavior, such as testosterone’s ability to predict 5  
 6 aggression (cortisol and serotonin do too).<sup>49</sup> Even some of connec- 6  
 7 tomics’ biggest proponents seem to see this problem, though perhaps 7  
 8 not the consequences of it: “The ability of pharmacological agents to 8  
 9 rapidly induce sleep, tranquility, excitement, hallucinations and so on 9  
 10 means that the behavioral state can be dramatically altered probably 10  
 11 without any modification to the connectome.”<sup>50</sup> Of course, to be ex- 11  
 12 cited or tranquil is to be in a particular psychological state. 12

13 This is not to say that the biological theory is true. In consciousness 13  
 14 studies—as elsewhere in science—ruling out false theories is the goal, 14  
 15 whereas finding true theories is a bit idealistic. Perhaps the best re- 15  
 16 sponse for functionalists is to become subneural functionalists, where 16  
 17 the properties that matter for functional realization are below the 17  
 18 neuronal level—perhaps far below (for example, perhaps biochem- 18  
 19 ical or subatomic). This would be an interesting discovery—the idea 19  
 20 that neural properties are not the functional realizers of the mind 20  
 21 is, at the very least, very surprising. Moreover, the resulting dilemma 21  
 22 itself, that one is forced to be a biological theorist or a subneural func- 22  
 23 tionalist, is an interesting-enough endpoint. 23

24 But becoming a subneural functionalist is also rather destructive 24  
 25 to the idea that WBE is the best chance to achieve superintelligence 25  
 26 or immortality. Subneural functionalism contravenes the “nonorgani- 26  
 27 cism” that allow futurists to champion WBE in the first place. The rea- 27  
 28 son WBE is so appealing to transhumanists, futurists, and the like is 28  
 29 that it seems much less farfetched than all the other routes to posthu- 29  
 30 man intelligence. WBE is supposed to be a data saver; it supposes that 30  
 31 all we need to do is upload the functional properties, so we do not 31  
 32 need to know how the whole brain works (or how the whole body 32  
 33 works, or how the whole species works, or how the whole universe 33  
 34 works). However, the lower the level of the functional properties, the 34  
 35 more we will need to know (and the more information we would need 35  
 36 to upload), meaning we would be much further away from achieving 36  
 37 uploading than even skeptics might assume. If the relevant level of de- 37  
 38 tail demands molecule-for-molecule duplication, then WBE looks to 38  
 39 39

40 <sup>48</sup> Morgan and Lichtman, “Why Not Connectomics?,” *op. cit.*, p. 496. 40

41 <sup>49</sup> E. Montoya et al., “Testosterone, Cortisol, and Serotonin as Key Regulators of So- 41  
 42 cial Aggression: A Review and Theoretical Perspective,” *Motivation and Emotion*, xxxvi, 42  
 1 (2012): 65–73.

43 <sup>50</sup> Morgan and Lichtman, “Why Not Connectomics?,” *op. cit.*, p. 497. 43

1 be entirely infeasible as an engineering project in even the medium- 1  
2 to-far term (and possibly computationally intractable). 2

3 So, if subneural functionalism is true, then the viability of WBE is 3  
4 in trouble. But we can go further still, for if the biological theory is 4  
5 true, much deeper theoretical revisions will be needed. If the biolog- 5  
6 ical theory is true, multiple realizability, computationalism, and even 6  
7 functionalism cannot be true of the entire mind. These theories may 7  
8 be true of propositional attitudes, or some other aspect of cognition, 8  
9 but they are not true of consciousness, in which case the mere pos- 9  
10 sibility of machine consciousness and WBE is imperiled. This moral 10  
11 has not been lost on the proponents of the biological theory, such 11  
12 as Block. He writes, “The biological theory says that only machines 12  
13 that have the right biology can have consciousness, and in that sense 13  
14 the biological account is less friendly to machine consciousness.”<sup>51</sup> Of 14  
15 course, we are not in a position to say that the biological theory is 15  
16 true. But it is enough to note that it is, at this time, still very much 16  
17 alive, and one of the very few live theories we have of consciousness, 17  
18 even if it is extremely underspecified. 18  
19

#### 20 IV. SHOULD WBE OPTIMISTS CARE ABOUT CONSCIOUSNESS? 20

21 In discussing the viability of WBE, Sandberg opined that “there 21  
22 doesn’t seem to be any convincing knock-down arguments within the 22  
23 philosophy of mind against WBE.”<sup>52</sup> Although there is not a knock- 23  
24 down argument against it, there is reason to have serious skepticism 24  
25 about WBE’s viability, and this, in turn, reveals some deeper problems 25  
26 in the metaphysics of mind. 26

27 Before concluding, let us take a step back to consider the big pic- 27  
28 ture: what did we want WBE for anyway? Only two ends have been put 28  
29 forth. The first is as a step toward achieving superintelligence, and the 29  
30 second is for achieving immortality. I take these in turn. 30

31 As a reminder, the route to superintelligence went through using 31  
32 WBE to upload human-level intelligence. Once we have a cheap and 32  
33 easy way to produce and store human intelligence, we can create an 33  
34 enormous amount of uploads and then put them to the task of dis- 34  
35 covering the breakthroughs that can lead to superintelligence. 35

36 How much would consciousness matter for this program? Say the 36  
37 biological theory is only true for phenomenal consciousness. Could 37  
38 the rest of cognition then be captured by the connectome, in which 38  
39 case WBE could still lead to superintelligence? The question turns, 39  
40

42 <sup>51</sup>Block, “Comparing the Major Theories of Consciousness,” *op. cit.*, p. 1119. 42

43 <sup>52</sup>Sandberg, “Feasibility of Whole Brain Emulation,” *op. cit.*, p. 261. 43

1 in part, on whether there can be intentionality without phenomenol- 1  
 2 ogy. Having some unconscious intentional states—like beliefs—is a 2  
 3 position that is held commonly enough.<sup>53</sup> But could there also moti- 3  
 4 vation, or desire, without any phenomenology? That seems much less 4  
 5 clear. What it is to desire something seems to involve feeling a cer- 5  
 6 tain way. Likewise, what it is to be motivated has an aversive quality 6  
 7 to it, which is just to say that some motivations appear to have some 7  
 8 phenomenology. 8

9 If we want uploads to do anything, they will have to be motivated.<sup>54</sup> 9  
 10 Cognition without conation is just a spinning wheel connected to 10  
 11 nothing. Having a billion more human-level intellects available to 11  
 12 work on a problem will only help solve the problem if they are de- 12  
 13 signed to solve the problem or motivated to do the work. Part of the 13  
 14 appeal of uploads is that we would not have to design any particu- 14  
 15 lar goal for them, for doing so takes us far beyond merely uploading 15  
 16 a connectome. Since we will not be able to design uploads with the 16  
 17 goal of solving any particular problem, uploads will only act if they 17  
 18 are intrinsically motivated to. If they have no motivations, then they 18  
 19 will not do anything on their own. 19

20 The problems for WBE get even worse. Many theories of the at- 20  
 21 titudes dictate that to have any beliefs at all, one must have other 21  
 22 propositional attitudes, particularly desires and motivations.<sup>55</sup> If there 22  
 23 are no desires, then uploads may not even have beliefs, for, so the 23  
 24 thought goes, part of the functional role that is constitutive of beliefs 24  
 25 is that they interact with desires to cause action. If uploads do not 25  
 26 have beliefs, it is hard to see how they could ever engage in thinking 26  
 27 as they would lack the premises of thoughts (and the desires to go 27  
 28 through the bother of transitioning from thought to thought).<sup>56</sup> 28

29 There is an even more exotic argument against the existence of be- 29  
 30 liefs that are totally disconnected from phenomenology. It starts by 30  
 31 noting that our beliefs matter to us. When we encounter disconfirm- 31  
 32 ing information it *hurts* and immediately causes us to readjust our 32  
 33 33

34 <sup>53</sup>Eric Mandelbaum, “Thinking Is Believing,” *Inquiry*, LVII, 1 (2014): 55–96. For a de- 34  
 35 fense of phenomenal intentionality see Uriah Kriegel, “The Phenomenal Intentionality 35  
 36 Research Program,” in Uriah Kriegel, ed., *Phenomenal Intentionality* (New York: Oxford 36  
 37 University Press, 2013), pp. 1–26. 37

38 <sup>54</sup>One may argue that cars and calculators do things without being motivated, but 38  
 39 they do so at the behest of intelligent, motivated designers and users. Even Bostrom’s 39  
 40 paperclip maximizer has to be seen as either having the motivation to turn everything 40  
 41 into paperclips or having been given the function to do so. Nick Bostrom, *Superintelli-* 40  
 42 *gence: Paths, Dangers, Strategies* (Oxford: Oxford University Press, 2014). 41

41 <sup>55</sup>Jerry A. Fodor, *Psychosemantics: The Problem of Meaning in the Philosophy of Mind* 41  
 42 (Cambridge, MA: MIT Press, 1987). 42

42 <sup>56</sup>Jake Quilty-Dunn and Eric Mandelbaum, “Inferential Transitions,” *Australasian* 42  
 43 *Journal of Philosophy*, xcvi, 3 (2018): 532–47. 43

1 beliefs, often perversely increasing credence in the proposition under  
2 attack.<sup>57</sup> Some hold that this is a defining feature of belief, so  
3 that any state that did not act this way would not be a belief.<sup>58</sup> If this is  
4 right, then if the connectome did not include valences, uploads could  
5 not have beliefs. And this argument generalizes for any mental state  
6 where valence plays a constitutive role.

7 If uploads lacked beliefs and desires, then they would just be gi-  
8 gant calculators that we neither knew how to control nor understood  
9 the mechanics of. Recall that the appeal of WBE was its nonorgani-  
10 cism, which allows that we could copy the brain without needing to  
11 understand how all of it works—this is what was supposed to move  
12 up the timetable of feasibility for WBE versus any other technologies.  
13 Then once we had the uploads, we could reason with them the way  
14 we would with any belief-/desire-based agent. But if uploads do not  
15 have the normal attitudes, we will have no idea how to motivate them  
16 to do anything—it is not even clear that they would be able to be moti-  
17 vated. In that case, we would have to go back to a more fine-grained  
18 stance to affect their behavior, which would demand another concep-  
19 tual breakthrough.

20 WBE’s promise for immortality raises even murkier questions. We  
21 generally think the issue of immortality and uploads boils down to  
22 the question of whether uploading your mind without consciousness  
23 would suffice for immortality. But even smaller questions about con-  
24 sciousness fester: might one’s particular type of phenomenology mat-  
25 ter for capturing identity? Does one’s character intimately involve the  
26 kind of phenomenology they have? Maybe you could be you even  
27 with a different character. This is not totally implausible—people can  
28 change their personality throughout their lifespan (though whether  
29 that actually makes a change in personhood is tendentious).<sup>59</sup> Yet  
30 some of the properties that seem deeply central to our self-conception  
31 would be left out. Above we noted that tranquility, excitement, and  
32 the like will be left out of the connectome. These properties are  
33 plainly not just properties at the edges of our identity but instead  
34 are often integral to who we are. People think of themselves as, for  
35 example, deeply energetic, or extremely calm and patient. But those  
36 personality traits would be left out of the connectome. Could your  
37

39 <sup>57</sup>Eric Mandelbaum, “Troubles with Bayesianism: An Introduction to the Psycholog-  
40 ical Immune System,” *Mind and Language*, xxxiv, 2 (2019): 141–57.

41 <sup>58</sup>Nicolas Porot and Eric Mandelbaum, “The Science of Belief: A Progress Report,”  
42 *Wiley Interdisciplinary Reviews: Cognitive Science*, xii, 2 (2021): e1539.

43 <sup>59</sup>Nina Strohminger and Shaun Nichols, “The Essential Moral Self,” *Cognition*,  
44 cxxxii, 1 (2014): 159–71.

1 connectome duplicate you even if, for example, it was a sickly sloth 1  
2 while you are a dynamo bursting at the seams with energy and ideas? 2

3 Even without taking a stand on what exactly personal identity 3  
4 amounts to, it appears that what it is like to be you does have some 4  
5 bearing on what it is to be you. And if that is the case, then the biggest 5  
6 roadblock to the grandiose promise of WBE—uploads—is that our 6  
7 biological machinery itself may be responsible for a good deal of our 7  
8 cognitive life. The problem is not just that, for example, you see deep 8  
9 purple whereas the upload version of you would experience periwinkle 9  
10 It is that to exist as you would involve some of the full panoply of 10  
11 emotions, feelings, depths, and depravities of everyday life, and these 11  
12 would be left out of the uploads. 12

13 This does not mean that we should endorse Mysterianism or be 13  
14 sure that uploading is necessarily impossible. The world never ceases 14  
15 to surprise. Perhaps one day we will be able to upload full wiring di- 15  
16 agrams into hardware just like ours. But if so, that would no longer 16  
17 be emulating whole brains, but cloning and recreating them from 17  
18 scratch, in which case the feasibility of achieving it should seem that 18  
19 much further off than current futurists prognosticate. 19

20 ERIC MANDELBAUM 20

21 City University of New York 21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43