

Integrating Philosophy of Understanding with the Cognitive Sciences

1 Kareem Khalifa^{1*}, Farhan Islam², J.P. Gamboa³, Daniel A. Wilkenfeld⁴, Daniel Kostić⁵

- ² ¹Department of Philosophy, Middlebury College, Middlebury, VT, USA
- ³ ²Independent Scholar, Madison, WI, USA
- ⁴ ³Department of History and Philosophy of Science, University of Pittsburgh, Pittsburgh, PA
- ⁵ ⁴Department of Acute and Tertiary Care, University of Pittsburgh School of Nursing, Pittsburgh, PA
- ⁶ ⁵Institute for Science in Society (ISiS), Radboud University, Nijmegen, The Netherlands
- 7

8 * Correspondence:

- 9 Kareem Khalifa
- 10 kkhalifa@middlebury.edu

Keywords: explanation, understanding, mechanism, computation, topology, dynamic systems, pluralism. (Min.5-Max. 8)

13 Abstract

14 We provide two programmatic frameworks for integrating philosophical research on understanding

15 with complementary work in computer science, psychology, and neuroscience. First, philosophical

- 16 theories of understanding have consequences about how agents should reason if they are to
- 17 understand that can then be evaluated empirically by their concordance with findings in scientific
- 18 studies of reasoning. Second, these studies use a multitude of explanations, and a philosophical
- 19 theory of understanding is well suited to integrating these explanations in illuminating ways.

20 **1** Introduction

21 Historically, before a discipline is recognized as a science, it is a branch of philosophy. Physicists

and chemists began their careers as "natural philosophers" during the Scientific Revolution. Biology

and psychology underwent similar transformations throughout the nineteenth and early twentieth

centuries. So, one might think philosophical discussions of understanding will be superseded by a

25 "science of understanding."

While we are no great forecasters of the future, we will suggest that philosophical accounts of understanding can make two important scientific contributions. First, they provide a useful repository hypotheses that can be operationalized and tested by scientists. Second, philosophical accounts of understanding can maximize templates for unificing a variety of accounts of understanding can maximize templates for unificing a variety of accounts of

- 29 understanding can provide templates for unifying a variety of scientific explanations.
- 30 We proceed as follows. Section 2 presents these two frameworks for integrating philosophical 31 ideas about understanding with scientific research. Sections 3 discusses the first of these frameworks,
- ideas about understanding with scientific research. Sections 3 discusses the first of these frameworks,
 in which philosophical theories of understanding propose hypotheses that are tested and refined by
- the cognitive sciences. Section 4 discusses the second, in which considerations of understanding

- 34 provide criteria for integrating different scientific explanations. Both of our proposals are intended to
- be programmatic. We hope that many of the relevant details will be developed in future work. 35

Two Frameworks for Integration 36 2

- 37 As several reviews attest (Baumberger, 2014; Baumberger, Beisbart, & Brun, 2016; Gordon, 2017;
- 38 Grimm, 2021; Hannon, 2021), understanding has become a lively topic of philosophical research
- 39 over the past two decades. While some work has been done to integrate these ideas with relevant
- findings from computer science, psychology, and neuroscience, these interdisciplinary pursuits are 40
- relatively nascent. While other frameworks are possible and should be developed, we propose two 41 ways of effecting a more thoroughgoing synthesis between philosophy and these sciences (Figure 1). 42
- 43
- In the first framework for integrating philosophy with the cognitive sciences—what we call 44 *naturalized epistemology of understanding* (Figure 1A)—the philosophy of understanding provides
- conjectures about reasoning that are tested and explained by the relevant sciences. In the second 45
- integrative framework—understanding-based integration (Figure 1B)—the philosophy of 46
- 47 understanding provides broad methodological guidelines about how different kinds of scientific
- explanation complement each other. The two proposals are independent of each other: those 48
- unpersuaded by one may still pursue the other. We discuss each in turn. 49

50 A. Naturalized Epistemology of Understanding



51

52 B. Understanding-Based Integration



53

Figure 1. Two Ways to Integrate Philosophical Work on Understanding with Relevant 54 55 Sciences.

56 3 Naturalized Epistemology of Understanding

57 In epistemology, naturalism is the position that philosophical analyses of knowledge, justification,

and kindred concepts should be intimately connected with empirical science. Different naturalists 58

59 specify this connection in different ways; see Rysiew (2020) for a review. Given that philosophical

60 interest in understanding has only recently achieved critical mass, the more specific research program

of a naturalized epistemology of understanding is nascent. We propose to organize much existing 61

62 work according to the framework in Figure 1A. More precisely, philosophical theories of

understanding propose how reasoning operates in understanding (Section 3.1), and these proposals 63

are constrained by explanations and empirical tests found in sciences that study this kind of reasoning 64

65 (Section 3.2).

66 3.1 Philosophical Theories Propose Reasoning in Understanding (I)

Two kinds of understanding have garnered significant philosophical attention: explanatory 67

68 understanding (Greco, 2013; Grimm, 2010, 2014; Hills, 2015; Khalifa, 2012, 2013a, 2013b, 2017;

Kuorikoski & Ylikoski, 2015; Potochnik, 2017; Strevens, 2013) and objectual understanding 69

70 (Baumberger, 2019; Baumberger & Brun, 2017; Carter & Gordon, 2014; Dellsén, 2020; Elgin, 2004,

2017; Kelp, 2015; Kvanvig, 2003; Wilkenfeld, 2021). Explanatory understanding involves 71

72 understanding why or how something is the case. (For terminological convenience, subsequent

references to "understanding-why" are elliptical for "understanding-why or -how".) Examples 73

include understanding why Caesar crossed the Rubicon and understanding how babies are made. 74

75 Objectual understanding is most easily recognized by its grammar: it is the word "understanding"

- followed immediately by a noun phrase, e.g., understanding Roman history or understanding human 76
- 77 reproduction. Depending on the author, the objects of objectual understanding are taken to be subject
- 78 matters, phenomena, and for some authors (e.g., Wilkenfeld, 2013), physical objects and human
- 79 behaviors. For instance, it is natural to think of Roman history as a subject matter but somewhat
- counterintuitive to think of it as a phenomenon. It is more natural to think of, e.g., the unemployment 80 81 rate in February 2021 as a phenomenon than as a subject matter. Human reproduction, by contrast,
- can be comfortably glossed as either a subject matter or a phenomenon. 82

Integrating Philosophy of Understanding

83 To clarify what they mean by explanatory and objectual understanding, philosophers have

84 disambiguated many other senses of the English word "understanding." Frequently, these senses are

briefly mentioned to avoid confusion but are not discussed at length. They are listed in Table 1.

- 86 Scientists may find these distinctions useful when characterizing the kind of understanding they are
- 87 studying. That said, we will focus on explanatory understanding hereafter. Thus, unless otherwise
- noted, all subsequent uses of "understanding" refer exclusively to explanatory understanding.

Kind of understanding	Typical Complement	Examples
Propositional	that + declarative sentence	I understand that you might not enjoy reading this book.
Broad Linguistic	name of a language	Schatzi understands German.
Narrow Linguistic	what + a linguistic expression + means	Schatzi understands what "Ich bin ein Berliner" means.
Procedural	how + infinitive	Miles understands how to play trumpet.
Non-explanatory Interrogative	embedded question that does not seek an explanation as its answer (most who, where, what, and when questions)	I understand who my friends are. I understand where my friends will be going. I understand what my friends are doing. I understand when my friends need a good laugh.

89

Table 1. Kinds of understanding that philosophers infrequently discuss (Khalifa 2017, 2)

90 Virtually all philosophers agree that one can possess an accurate explanation without 91 understanding it, e.g., through rote memorization. In cases such as this, philosophers widely agree that the lack of understanding is due to the absence of significant inferential or reasoning abilities. 92 93 However, philosophers disagree about which inferences characterize understanding. Three broad 94 kinds of reasoning have emerged. First, some focus on the reasoning required to construct or consider explanatory models (De Regt, 2017; Newman, 2012, 2013, 2015). Second, others focus on 95 96 the reasoning required to evaluate those explanatory models (Khalifa, 2017). On both these views, explanatory models serve as the conclusions of the relevant inferences. However, the third and most 97 prominent kind of reasoning discussed takes explanatory information as *premises* of the relevant 98 99 reasoning-paradigmatically the inferences about how counterfactual changes in the explanatory variable or explanans would result in changes to the dependent variable or explanandum ((Bokulich, 100 2011; Grimm, 2010, 2014; Hills, 2015; Hitchcock & Woodward, 2003; Kuorikoski & Ylikoski, 101 102 2015; Le Bihan, 2016; Potochnik, 2017; Rice, 2015; Verreault-Julien, 2017; Wilkenfeld, 2013;

103 Woodward, 2003). This is frequently referred to as the ability to answer "what-if-things-had-been-

104 different questions." Many of these authors discuss all three of these kinds of reasoning—which we

105 call explanatory consideration, explanatory evaluation, and counterfactual reasoning—often without

106 explicitly distinguishing them in the ways we have here.

107 **3.2** Scientific Studies of Reasoning's Contributions to the Philosophy of Understanding (II)

108 A naturalized epistemology of understanding begins with the recognition that philosophers do not

- 109 have a monopoly on studying these kinds of reasoning. Computer scientists, psychologists, and
- 110 neuroscientists take explanatory and counterfactual reasoning to be important topics of research.
- 111 Undoubtedly, each discipline has important insights and contributions. Moreover, these scientific
- 112 disciplines may raise interesting questions about understanding that are not on the current
- 113 philosophical agenda.

114 Cognitive psychological investigations into the nature of explanation and understanding 115 frequently focus on the role of those states in our cognitive lives. To the extent that one can derive a 116 general lesson from this literature, it is probably that both having and seeking explanations aid other 117 crucial cognitive tasks such as prediction, control, and categorization. Developmental psychologists 118 argue that having proper explanations promotes survival, and that at least the sense of understanding 119 evolved to give us an immediate reward for gaining such abilities (Gopnik, 1998). In cognitive 120 psychology, Koslowski, Marasia, Chelenza, and Dublin (2008) have argued that having an 121 explanation better enables thinkers to incorporate evidence into a causal framework. Lombrozo and 122 collaborators have done extensive empirical work investigating the epistemic advantages and 123 occasional disadvantages of simply being prompted to explain new data. They find that under most 124 normal circumstances trying to seek explanations enables finding richer and more useful patterns 125 (Williams & Lombrozo, 2010). This work also has the interesting implication that the value of 126 explanation and understanding depends on the extent to which there are genuine patterns in the 127 world, with fully patterned worlds granting the most advantages from prompts to explain (ibid.), and 128 more exception-laden worlds providing differential benefits (Kon & Lombrozo, 2019). It has also 129 been demonstrated that attempts to explain can (perhaps counterintuitively) systematically mislead. 130 For example, attempts to explain can lead to miscategorization and inaccurate predictions when there 131 are no real patterns in the data (Williams, Lombrozo, & Rehder, 2013). Similarly, laypeople can be 132 misguided by the appearance of irrelevant neuroscientific or otherwise reductive explanations 133 (Hopkins, Weisberg, & Taylor, 2016; Weisberg, Keil, Goodstein, Rawson, & Gray, 2008). In more 134 theoretical work, Lombrozo (2006; Lombrozo & Wilkenfeld, 2019) considers how different kinds of 135 explanation can lead to understanding that is either more or less tied to specific causal pathways connecting explananda and explanantia versus understanding focused on how different pathways can 136 137 lead to the same end result. Thagard (2012) has argued that explanatory reasoning is key to science's 138 goals both intrinsically and as they contribute to truth and education.

139 One recent thread in the cognitive science and philosophy of understanding combines insights 140 from information theory and computer science to characterize understanding in terms of data 141 compression. Data compression (Grünwald, 2004) involves the ability to produce large amounts of 142 information from relatively shorter hypotheses and explicitly encoded data sets-in computer science 143 and model-centric physics, there is a burgeoning sense that understanding is tied to pattern 144 recognition and data compression. Petersen (ms) helpfully documents an array of such instances. Li 145 and Vitányi (2008) use compression and explanation almost interchangeably, and at some points 146 even suggest a possible equivalence between compression and the scientific endeavor generally, as in 147 Davies (1990). Tegmark (2014) likewise connects the notion of compression with the explanatory

148 goals of science. Wilkenfeld (2019) translates the importance of compression to good scientific (and

149 non-scientific) understanding into the idiom of contemporary philosophy of science. While part of

150 the inspiration characterizing understanding in terms of compression comes from the traditional 151 "unificationist" philosophical position that understanding involves having to know fewer brute facts

151 (Friedman, 1974) or argument patterns (Kitcher, 1989), the introduction of compression helps evade

some objections to unificationist views, such as the fact that such views require explanations to be

arguments (Woodward, 2003) and the fact that they allow for understanding via unification that no

actual human agent can readily use (Humphreys, 1993). (Compression as a marker for intelligence

156 has come under recent criticism (e.g., Chollet, 2019) as only accounting for past data and not future

uncertainties; we believe Wilkenfeld's (2019) account evades this criticism by defining the relevant

158 compression partially in terms of usefulness, but defending that claim is beyond the scope of this 159 paper).

160 There has also been more direct work on leveraging insights from computer science in order 161 to try to build explanatory schema and even to utilize those tools to reach conclusions about true 162 explanations. Schank (1986) built a model of computerized explanations in terms of scripts and 163 designed programs to look for the best explanations. Similarly, Thagard (1989, 1992, 2012)—who 164 had previously (1978) done seminal philosophical work on good-making features of explanation and 165 how they should guide theory choice—attempted to automate how computers could use 166 considerations of explanatory coherence to make inferences about what actually occurred.

167 One underexplored area in the philosophy of understanding and computer science is the 168 extent to which neural nets and deep learning machines can be taken to understand anything. While 169 Turing (1950) famously argued that a machine that could behave sufficiently close to a person could 170 thereby think (and thus, perhaps, understand), many argue that learning algorithms are concerned 171 with prediction as opposed to understanding. The most extreme version of this position is Searle's 172 (1980) claim that computers by their nature cannot achieve understanding because it requires 173 semantic capacities when manipulating symbols (i.e., an ability to interpret symbols and operations, 174 and to make further inferences based on those interpretations). Computers at best have merely 175 syntactic capabilities (they can manipulate symbols using sets of instructions, without understanding 176 the meaning of either symbols or operation upon them). However, at the point where deep learning machines have hidden representations (Korb, 2004), can generate new (seemingly theoretical) 177 178 variables (ibid.), and can be trained to do virtually any task to which computer scientists have set 179 their collective minds (including what looks from the outside like abstract reasoning in IBM's 180 Watson and their Project Debater), it raises vital philosophical questions regarding on what basis we

181 can continue to deny deep learning machines the appellation of "understander".

182 Elsewhere in cognitive science, early psychological studies of reasoning throughout the 1960s 183 and 1970s focused on deductive reasoning and hypothesis testing (Osman, 2014). A major influence 184 on this trajectory was Jean Piaget's (1952) theory of development, according to which children develop the capacity for hypothetico-deductive reasoning around age 12. The kinds of reasoning 185 186 studied by psychologists then expanded beyond their logical roots to include more humanistic categories such as moral reasoning (Kohlberg, 1958). The psychology literature offers a rich body of 187 188 evidence demonstrating how people reason under various conditions. For example, there is ample 189 evidence that performance on reasoning tasks is sensitive to the semantic content of the problem 190 being solved. One interpretation of this phenomenon is that in some contexts, people do not reason 191 by applying content-free inference rules (Cheng & Holyoak, 1985; Cheng, Holyoak, Nisbett, &

192 Oliver, 1986; Holyoak & Cheng, 1995). This empirical possibility is of particular interest for 193 philosophers. In virtue of their (sometimes extensive) training in formal logic, philosophers' 194 reasoning practices may be atypical of the broader population. This in turn may bias their intuitions about how "people" or "we" reason in various situations, including when understanding. Another 195 196 issue raised by sensitivity to semantic content is how reasoning shifts depending on the object of 197 understanding. Although the distinctions explicated by philosophers (e.g., explanatory vs. objectual 198 understanding) are clear enough, it is an open empirical question whether and how reasoning differs 199 within these categories depending on the particular object and other contextual factors. As a final example, a further insight from psychology is that people may have multiple modes of reasoning that 200 can be applied to the very same problem. Since Wason and Evans (1974) suggested the idea, dual-201 202 process theories have dominated the psychology of reasoning.¹ Although both terminology and 203 precise hypotheses vary significantly among dual-process theories (Evans, 2011, 2012), the basic 204 idea is that one system of reasoning is fast and intuitive, relying on prior knowledge, while another is slow and more cognitively demanding. Supposing two or more systems of reasoning can be deployed 205 206 in the same situation, one important consideration is how they figure in theories about the reasoning 207 involved in understanding. To the extent that philosophical accounts are not merely normative but 208 also aim at describing how people actually reason when understanding, psychological studies provide 209 valuable empirical constraints and theoretical considerations.

210 With the aid of techniques for imaging brains while subjects perform cognitive tasks, 211 neuroscientists have also made great progress in recent decades on identifying regions of the brain 212 involved in reasoning. While that is certainly a worthwhile goal, it may seem tangential to 213 determining the kind of reasoning that characterizes understanding. Here, we suggest two ways in 214 which findings from neuroscience may help with this endeavor. First, neuroscientific evidence can 215 help resolve debates where behavioral data underdetermine which psychological theory is most 216 plausible. More precisely, in cases where competing psychological models of reasoning make the 217 same behavioral predictions, they can be further distinguished by the kinds of neural networks that 218 would implement the processes they hypothesize (Operskalski & Barbey, 2017). For example, Goel, 219 Buchel, Frith, and Dolan (2000) designed a functional magnetic resonance imaging (fMRI) 220 experiment to test the predictions of dual mechanism theory vs. mental model theory. According to the former, people have distinct mechanisms for form- and content-based reasoning, and the latter 221 222 should recruit language processing structures in the left hemisphere. Mental model theory, by 223 contrast, claims that reasoning essentially involves iconic representations, i.e., non-linguistic 224 representations whose structure corresponds to the structure of whatever they represent (Johnson-225 Laird, 2010). In early formulations of the theory, it was assumed that different kinds of reasoning 226 problems depend on the same visuo-spatial mechanisms in the right hemisphere (Johnson-Laird, 227 1995). Goel et al. (2000) tested the theories against one another by giving subjects logically 228 equivalent syllogisms with and without semantic content. As expected, behavioral performance was 229 similar in both conditions. Neither theory predicts significant behavioral differences. Consistent with 230 both theories, the content-free syllogisms engaged spatial processing regions in the right hemisphere. However, syllogisms with semantic content activated a left hemisphere ventral network that includes 231

¹ Though see Keren and Schul (2009), Osman (2004), and Stephens, Dunn, and Hayes (2018) for ⁷ examples of criticisms.

232 language processing structures like Broca's area. Unsurprisingly, proponents of mental models have

- disputed the interpretation of the data (Kroger, Nystrom, Cohen, & Johnson-Laird, 2008). We do not
- take a stance on the issue here. We simply raise the case because it illustrates how neuroscience can
- contribute to debates between theories of reasoning pitched at the psychological level.

236 Neuroscientific evidence can also guide the revision of psychological models of 237 understanding and reasoning. The broader point is about cognitive ontology. In the sense we mean here, a cognitive ontology is a set of standardized terms which refer to the entities postulated by a 238 239 cognitive theory (Janssen, Klein, & Slors, 2017). The point of developing a cognitive ontology is to 240 represent the structure of psychological processes and facilitate communication through a shared 241 taxonomy. One role for neuroscience is to inform the construction of cognitive ontologies. Price and 242 Friston (2005), for instance, defend a strong bottom-up approach. In their view, components in a cognitive model (e.g., a model of counterfactual reasoning) should be included or eliminated 243 244 depending on our knowledge of functional neuroanatomy. Others agree that neuroscience has a 245 crucial role to play in theorizing about cognitive architecture but reject that it has any special authority in this undertaking (Poldrack & Yarkoni, 2016; Sullivan, 2017). We take no position here 246 on how exactly neuroscience should influence the construction of cognitive models and ontologies. 247 248 Instead, we highlight this important interdisciplinary issue to motivate the potential value of 249 neuroscience for models of understanding and the reasoning involved in it, including those developed

250 by philosophers.

251 4 Philosophical Theories of Understanding Integrate Scientific Explanations (III)

252 Thus, there appear to be ample resources for a naturalized epistemology of understanding, in which

explanations and empirical tests from the cognitive sciences empirically constrain philosophical proposals about the kinds of reasoning involved in understanding. However, we offer a second and

proposals about the kinds of reasoning involved in understanding. However, we offer a second and distinct proposal for how the philosophy of understanding can inform scientific practice: as an

256 account of how different explanations can be integrated (Figure 1B).

257 Such integration is needed when different explanations of a single phenomenon use markedly

- 258 different vocabularies and concepts. This diversity of explanations is prevalent in several sciences—
- 259 including the cognitive sciences. To that end, Section 4.1 presents different kinds of explanations
- 260 frequently found in the cognitive sciences. Whether these different explanations are complements or
- 261 competitors to each other raises several issues that are simultaneously methodological and
- 262 philosophical. To address these issues, Section 4.2 presents a novel account of explanatory
- 263 integration predicated on the idea that explanations are integrated to the extent that they collectively 264 promote understanding. To illustrate the uniqueness of this account, Section 4.3 contrasts our account
- 264 promote understanding. To industrate the uniqueness of this account, Section 4.5 contra 265 of integration with a prominent alternative in the philosophical literature.
- Before proceeding, two caveats are in order. First, although we focus on the cognitive sciences, the account of explanatory integration proposed here is perfectly general. In principle, the same account could be used in domains ranging from particle physics to cultural anthropology. Second, our aim is simply to show that our account of integration enjoys some initial plausibility; a more thoroughgoing defense exceeds the current paper's scope.

271 **4.1 A Variety of Scientific Explanations**

- 272 Puzzles about explanatory integration arise only if there are explanations in need of integration, i.e.,
- explanations whose fit with each other is not immediately obvious. In this section, we provide
- examples of four kinds of explanations found in the cognitive sciences: mechanistic, computational,
- topological, and dynamical.

276 4.1.1 Mechanistic Explanations

277 Mechanistic explanations are widespread in the cognitive sciences (Bechtel & Richardson, 1993;

- 278 Craver, 2007; Craver & Tabery, 2019; Glennan, 2017; Illari & Williamson, 2010; Machamer,
- 279 Darden, & Craver, 2000). Despite extensive discussion in the philosophical literature, there is no
- consensus on the proper characterization of mechanisms or how exactly they figure in mechanistic
- 281 explanations.² For our purposes, we illustrate basic features of mechanistic explanations by focusing
- on Glennan's (2017, p. 17) minimal conception of mechanisms:
- A mechanism for a phenomenon consists of entities (or parts) whose activities and interactions are organized so as to be responsible for the phenomenon.
- 285 This intentionally broad proposal captures a widely held consensus among philosophers about
- conditions that are necessary for something to be a mechanism. Where they disagree is about further

details, such as the nature and role of causation, regularities, and levels of analysis involved in

- 288 mechanisms. At minimum, mechanistic explanations account for the phenomenon to be explained
- 289 (the *explanandum*) by identifying the organized entities, activities, and interactions responsible for it.
- 290 Consider the case of the action potential. A mechanistic explanation of this phenomenon specifies
- 291 parts such as voltage-gated sodium and potassium channels. It describes how activities of the parts,
- 292 like influx and efflux of ions through the channels, underlie the rapid changes in membrane potential.
- 293 It shows how these activities are organized such that they are responsible for the characteristic phases
- of action potentials. For example, the fact that depolarization precedes hyperpolarization is explained
- in part by the fact that sodium channels open faster than potassium channels. In short, mechanistic
- 296 explanations spell out the relevant physical details.
- Importantly, not all theoretical achievements in neuroscience are mechanistic explanations. As a
- 298 point of contrast, compare Hodgkin and Huxley's (1952) groundbreaking model of the action
- potential. With their mathematical model worked out, they were able to predict properties of action
- 300 potentials and neatly summarize empirical data from their voltage clamp experiments. However, as
- Hodgkin and Huxley (1952) explicitly pointed out, their equations lacked a physical basis. There is
- 302 some disagreement among philosophers about how we should interpret the explanatory merits of the
- 303 model (Craver & Kaplan, 2020; Favela, 2020a; Levy, 2014), but what is clear is that the Hodgkin and
- Huxley model is a major achievement that is *not* a mechanistic explanation of the action potential.
- We return to issues such as these in Section 4.3.

306 4.1.2 Computational Explanations

Mechanistic explanations are sometimes contrasted with other kinds of explanation. In the
 philosophical literature, computational explanations are perhaps the most prominent alternative.
 Computational explanations are frequently considered a subset of *functional explanations*. The latter

 $^{^{2}}$ See Craver (2014) for an overview of the latter issue.

310 explain phenomena by appealing to their function and the functional organization of their parts

- 311 (Cummins, 1975, 1983, 2000; Fodor, 1968). Insofar as computational explanations are distinct from
- 312 other kinds of functional explanations, it is because the functions to which they appeal involve
- 313 information processing. Hereafter, we focus on computational explanations.

314 In computational explanations, a phenomenon is explained in terms of a system performing a 315 computation. A computation involves the processing of input information according to a series of specified operations that results in output information. While many computational explanations 316 317 describe the object of computation as having representational content, some challenge this as a 318 universal constraint on computational explanations (Dewhurst, 2018; Fresco & Miłkowski, 2021; 319 Piccinini, 2015). We will use "information" broadly, such that we remain silent on this issue. Here, 320 "operations" refer to logical or mathematical manipulations on information such as addition, 321 subtraction, equation (setting a value equal to something), "AND", etc. For example, calculating n! 322 involves taking in input *n* and calculating the product of all natural numbers less than or equal to *n* and then outputting said product. Thus, we can explain why pressing "5", "!", "=", in sequence on a 323 calculator results in the display reading "120"; the calculator *computes* the factorial. 324

325 More detailed computational explanations of this procedure are possible. For example, the 326 calculator performs this computation by storing n and iteratively multiplying the stored variable by 327 one less than the previous iteration from n to 1. In this case, the operations being used are equation, 328 multiplication, and subtraction. The information upon which those operations are being performed 329 are the inputted value for n and the stored variable for the value of the factorial at that iteration.

330 4.1.3 Topological Explanations

331 In topological or "network" explanations, a phenomenon is explained by appeal to graph-332 theoretic properties. Scientists infer a network's structure from data, and then apply various graph-333 theoretic algorithms to measure its topological properties. For instance, clustering coefficients 334 measure degrees of interconnectedness among nodes in the same neighborhood. Here, a node's 335 *neighborhood* is defined as the set of nodes to which it is directly connected. An individual node's 336 local clustering coefficient is the proportion of edges within its neighborhood divided by the number of edges that could possibly exist between the members of its neighborhood. By contrast, a network's 337 338 global clustering coefficient is the ratio of closed triplets to the total number of triplets in a graph. A 339 triplet of nodes is any three nodes that are connected by at least two edges. An open triplet is 340 connected by exactly two edges; a *closed* triplet, by three. Another topological property, average (or 341 "characteristic") path length, measures the mean number of edges needed to connect any two nodes 342 in the network.

In their seminal paper, Watts and Strogatz (1998) applied these concepts to a family of graphs and showed how a network's topological structure determines its dynamics. First, *regular graphs* have both high global clustering coefficients and high average path length. By contrast, *random graphs* have low global clustering coefficients and low average path length. Finally, they introduced a third type of *small-world graph* with high clustering coefficient but low average path length.

Highlighting differences between these three types of graphs yields a powerful explanatory strategy. For example, because regular networks have larger average path lengths than small-world networks, things will "diffuse" throughout the former more slowly than the latter, largely due to the greater number of edges to be traversed. Similarly, because random networks have smaller clustering coefficients than small-world networks, things will also spread throughout the former more slowly than the latter, largely due to sparse interconnections within neighborhoods of nodes. Hence, *ceteris*

Integrating Philosophy of Understanding

354 *paribus*, propagation/diffusion is faster in small-world networks. This is because the fewer long-

355 range connections between highly interconnected neighborhoods of nodes shorten the distance

between neighborhoods of nodes that are otherwise very distant and enables them to behave as if they were first neighbors. For example, Watts and Strogatz showed that the nervous system of *C. elegans*

is a small-world network, and subsequent researchers argued that this system's small-world topology

- explains its relatively efficient information propagation (Bullmore & Sporns, 2012; Latora &
- 360 Marchiori, 2001).

361 **4.1.4 Dynamical Explanations**

In dynamical explanations, phenomena are accounted for using the resources of dynamic
 systems theory. At root, a system is dynamical if its state space can be described using differential
 equations, paradigmatically of the following form:

$$\dot{x}(t) = f(x(t); p, t)$$

Here, x is a vector (often describing the position of the system of interest), f is a function, t is time, and p is a fixed parameter. Thus, the equation describes the evolution of a system over time. In

dynamical explanations, these equations are used to show how values of a quantity at a given time.

and place would uniquely determine the phenomenon of interest, which is typically treated as values

370 of the same quantity at a subsequent time.

For example, consider dynamical explanations of why bimanual coordination—defined roughly as wagging the index fingers of both hands at the same time—is done either in- or anti-phase. Haken, Kelso, and Bunz (1985) use the following differential equation to model this phenomenon:

$$\frac{d\phi}{dt} = -asin\phi - 2bsin2\phi$$

Here ϕ is relative phase, having a value of either 0 degrees or 180 degrees (representing in- and antiphase conditions respectively) and b/a is the coupling ratio inversely related to the oscillations' frequency. The explanation rests on the fact that only the in- and anti-phase oscillations of the index fingers are basing of attraction

378 fingers are basins of attraction.

379 4.2 Understanding-Based Integration

380 Thus far, we have surveyed four different kinds of explanation—mechanistic, computational,

topological, and dynamical. Moreover, each seems to have some explanatory power for some

382 phenomena. This raises the question as to how these seemingly disparate kinds of explanation can be

383 integrated. We propose a new account of "understanding-based integration" (UBI) to answer this

question. A clear account of understanding is needed if it is to integrate explanations. To that end,

385 Section 4.2.1 presents Khalifa's (2017) model of understanding. Section 4.2.2 then extends this

account of understanding to provide a framework for explanatory integration.

387 4.2.1 An Account of Understanding

388 We highlight two reasons to think that Khalifa's account of understanding is especially promising as

a basis for explanatory integration. First, as Khalifa (2019) argues, his is among the most demanding

390 philosophical accounts of understanding. Consequently, it serves as a useful ideal to which scientists

391 should aspire. Second, this ideal is not utopian. This is especially clear with Khalifa's requirement

- that scientists evaluate their explanations relative to the best available methods and evidence. Indeed,
- 393 among philosophical accounts of understanding, Khalifa's account is uniquely sensitive to the

- 394 centrality of hypothesis testing and experimental design in advancing scientific understanding
- 395 (Khalifa, 2017, forthcoming), and thus makes contact with workaday scientific practices. In this
- 396 section, we present its three core principles.
- 397 Khalifa's first central principle is the *Explanatory Floor*:
- 398 Understanding why *Y* requires possession of a correct explanation of why *Y*.

399 The Explanatory Floor's underlying intuition is simple. It seems odd to understand why *Y* while

400 lacking a correct answer to the question, "Why Y?" For instance, the person who lacks a correct

401 answer to the question "Why do apples fall from trees?" doesn't understand why apples fall from

402 trees. Since explanations are answers to why-questions, the Explanatory Floor appears platitudinous.

403 Section 4.3.2.1 provides further details about correct explanation.

- The Explanatory Floor is only one of three principles comprising Khalifa's account and
 imposes only a necessary condition on understanding. By contrast, the second principle, the *Nexus Principle*, describes how understanding can improve:
- 407Understanding why Y improves in proportion to the amount of correct explanatory408information about Y (= Y's explanatory nexus) in one's possession.

409 To motivate the Nexus Principle, suppose that one person can correctly identify two causes of a fire,

410 and another person can only identify one of those causes. *Ceteris paribus*, the former understands

411 why the fire occurred better than the latter. Crucially in what follows, however, "correct explanatory

412 information" is not limited to correct explanations. The explanatory nexus also includes the

- 413 *relationships* between correct explanations. We return to these "inter-explanatory relationships"
- 414 below.

Furthermore, recall our earlier remark that gaps in understanding arise when one simply has
an accurate representation of an explanation (or explanatory nexus) without significant cognitive
ability. This leads to the last principle, the *Scientific Knowledge Principle*:

418 Understanding why *Y* improves as one's possession of explanatory information about *Y* bears
419 greater resemblance to scientific knowledge of *Y*'s explanatory nexus.

420 Once again, we may motivate this with a simple example. Consider two agents who possess the same 421 explanatory information that nevertheless differ in understanding because of their abilities to relate

421 explanatory information that neverticess differ in understanding because of their abilities to relate 422 that information to relevant theories, models, methods, and observations. The Scientific Knowledge

423 Principle is intended to capture this idea. Khalifa provides a detailed account of scientific knowledge

- 424 of an explanation:
- 425 An agent *S* has scientific knowledge of how/why *Y* if and only if there is some *X* such that *S*'s
 426 belief that *X* explains *Y* is the safe result of *S*'s scientific explanatory evaluation (SEEing).

427 The core notions here are safety and SEEing. Safety is an epistemological concept that requires an

428 agent's belief to not easily have been false given the way in which it was formed (Pritchard, 2009).

429 SEEing then describes the way a belief in an explanation should be formed to promote

430 understanding. SEEing consists of three phases:

431 1. *Considering* plausible potential explanations of how/why *Y*;

432 2. *Comparing* those explanations using the best available methods and evidence; and

- 433 3. Undertaking *commitments* to these explanations on the basis these comparisons.
- 434 Paradigmatically, commitment entails that one believes only those plausible potential
 435 explanations that are decisive "winners" at the phase of comparison.

Thus, scientific knowledge of an explanation is achieved when one's commitment to an explanation
could not easily have been false given the way that one considered and compared that explanation to
plausible alternative explanations of the same phenomenon.

439 **4.2.2 Understanding-Based Integration**

With our account of understanding in hand, we now argue that it provides a fruitful account of how different explanations, such as the ones discussed in Section 4.1, can be integrated. The Nexus Principle is the key engine of integration. As noted above, this principle states that understanding improves in proportion to the amount of explanatory information possessed. In the cognitive sciences, a multitude of factors explain a single phenomenon. According to the Nexus Principle, understanding improves not only when more of these factors are identified, but when the "interexplanatory relationships" between these factors are also identified.

447 One "inter-explanatory relationship" is that of *relative goodness*. Some explanations are better than others, even if both are correct. For instance, the presence of oxygen is explanatorily 448 449 relevant to any fire's occurrence. However, oxygen is rarely judged as the best explanation of a fire. Per the Nexus Principle, grasping facts such as these enhances one's understanding. Parallel points 450 451 apply in the cognitive sciences. For example, it has been observed that mental simulations that 452 involve episodic memory engage the default network significantly more than mental simulations that 453 involve semantic memory (Parikh, Ruzic, Stewart, Spreng, & De Brigard, 2018). Hence, episodic 454 memory better explains cases in which the default network was more active during a mental 455 simulation than does semantic memory.

456 However, correct explanations can stand in other relations than superiority and inferiority. 457 There are also "structural" relationships between different correct explanations. For instance, the 458 aforementioned explanation involving the default network contributes to a more encompassing 459 computational explanation of counterfactual reasoning involving three core stages of counterfactual 460 thought (Van Hoeck, Watson, & Barbey, 2015). First, alternative possibilities to the actual course of 461 events are mentally simulated. Second, consequences are inferred from these simulations. Third, 462 adaptive behavior and learning geared toward future planning and problem-solving occurs. The default network figures prominently in the explanation of (at least) the first of these processes (Figure 463 464 2).

As this example illustrates, grasping the relationships between different kinds of explanations can advance scientists' understanding. In Figure 2, a computational account of mental simulation explains certain aspects of counterfactual reasoning, but mental simulation is then explained mechanistically: the default network consists of parts (e.g., ventral medial prefrontal cortex, posterior cingulate cortex) whose activities and interactions (anatomical connections) are organized so as to be responsible for various phenomena related to mental simulations. Quite plausibly, scientific understanding increases when the relationship between these two explanations is discovered.

472 Importantly, this is but an instance of an indefinite number of other structures consisting of 473 inter-explanatory relationships (see Figure 3 for examples). In all of these structures, we assume that 474 for all *i*, X_i is a correct explanation of its respective explanatory. Intuitively, a person who could not

Integrating Philosophy of Understanding

- 475 distinguish these different explanatory structures would not understand *Y* as well as someone who
- 476 did. For instance, a person who knew that X_1 only explains Y through X_2 in Figure 3A, or that X_1 and
- 477 X_2 are independent of each other in Figure 3B, or that X_3 is a common explanation or "deep
- 478 determinant" of both X_1 and X_2 in Figure 3D, etc. seems to have a better understanding than a person
- who did not grasp these relationships. Undoubtedly, explanations can stand in other relationships thatfigure in the nexus.

481 Thus, the Nexus Principle provides useful guidelines for how different kinds of explanations 482 should be integrated. Moreover, we have already seen that different kinds of explanations can stand in fruitful inter-explanatory relationships, and that these relationships enhance our understanding. In 483 484 some cases, we may find that one and the same phenomenon is explained both mechanistically and 485 non-mechanistically, but one of these explanations will be better than another. As noted above, "better than" and "worse than" are also inter-explanatory relationships. So, the Nexus Principle 486 487 implies that knowing the relative strengths and weaknesses of different explanations enhances 488 understanding.

The Scientific Knowledge Principle also plays a role in UBI. Suppose that X_1 and X_2 are competing explanations of Y. SEEing would largely be achieved when, through empirical testing, X_1 was found to explain significantly more of Y's variance than X_2 . This gives scientists grounds for thinking X_1 better explains Y than X_2 and thereby bolsters our understanding of Y. Importantly, SEEing is also how scientists discover other inter-explanatory relationships. An example is the aforementioned study that identified the inter-explanatory relationships between episodic memory,

495 semantic memory, the default network, and mental simulation (Parikh et al., 2018).



496

497 Figure 2. Computational and Mechanistic Explanations Involved in Counterfactual Reasoning

498 Mental simulation (gray box) both contributes to the computational explanation of counterfactual

499 reasoning (black box) and is mechanistically explained by the activation of the default network.



500

501 Figure 3. Different Inter-Explanatory Relationships.

Letters at the head of an arrow denote phenomena to be explained; those at the tail, factors that do the explaining. For example, X_1 explains X_2 and X_2 explains Y in Figure 3A.

504 4.3 Mechanism-Based Integration

505 Aside from UBI, several other philosophical accounts of explanatory integration in the cognitive sciences are available (Kaplan, 2017; Miłkowski & Hohol, 2020). We provide some preliminary 506 507 comparisons with the most prominent of these accounts, which we call *mechanism-based integration* (MBI). According to strong MBI, all models in the cognitive sciences are explanatory only insofar as 508 509 they provide information about mechanistic explanations. In response, several critics of MBI-whom 510 we call *pluralists*—have provided examples of putatively non-mechanistic explanation (see Table 2). 511 When presented with putatively non-mechanistic explanations, e.g., of the computational, 512 topological, and dynamical varieties, mechanists (i.e., MBI's proponents) have two strategies 513 available. First, the negative strategy argues that closer scrutiny of the relevant sciences reveals the 514 putatively non-mechanistic explanation to be no explanation at all (Kaplan, 2011; Kaplan & Craver, 515 2011). The assimilation strategy argues that closer analysis of the relevant sciences reveals the putatively non-mechanistic explanation to be a mechanistic explanation, often of an elliptical nature 516 (Hochstein, 2016; Miłkowski, 2013; Piccinini, 2006, 2015; Piccinini & Craver, 2011; Povich, 2015; 517 518 Zednik, 2011). Mechanists inclined toward strong MBI frequently use the negative and assimilation 519 strategies in a divide-and-conquer-like manner: the negative strategy applies to some putatively non-520 mechanistic explanations and the assimilation strategy applies to the rest. However, more prevalent is 521 a *modest* form of MBI that simply applies these strategies to *some* putatively non-mechanistic explanations. 522

523 Modest MBI diverges from pluralism on a case-by-case basis. Such cases consist of an 524 explanation where the negative or assimilation strategy seems apt but stands in tension with other 525 considerations that suggest the model is both explanatory and non-mechanistic. On this latter front,

- 526 several pluralists argue that computational, topological, and dynamical explanations' formal and
- 527 mathematical properties are not merely abstract representations of mechanisms (Chirimuuta, 2018;
- 528 Darrason, 2018; Egan, 2017; Huneman, 2018; Lange, 2017; Rusanen & Lappi, 2016; Serban, 2015;
- van Rooij & Baggio, 2021; Weiskopf, 2011). Others argue that these explanations cannot (Chemero,
- 530 2009; Rathkopf, 2018; Silberstein & Chemero, 2013) or need not (Shapiro, 2019) be decomposed
- into mechanistic components or that they cannot be intervened upon in the same way that
- 532 mechanisms are intervened upon (Meyer, 2018; Ross, 2020; Woodward, 2013). Some argue that
- 533 these putatively non-mechanistic explanations are non-mechanistic because they apply to several 534 different kinds of systems that have markedly different mechanistic structures (Chirimuuta, 2014;
- different kinds of systems that have markedly different mechanistic structures (Chirimuuta, 2014;
 Ross, 2015). Pluralist challenges specific to different kinds of explanations can also be found (e.g.,
- 555 Koss, 2015). Pluralist challenges specific to different kinds of explanations can also be found (
- 536 Kostić, 2018; Kostić & Khalifa, manuscript).

537 In what follows, we will show how UBI is deserving of further consideration because it
538 suggests several plausible alternatives to the assimilation and negative strategies. As such, it contrasts
539 with both strong and modest MBI. While we are partial to pluralism, our discussion here is only

540 meant to point to different ways in which mechanists and pluralists can explore the issues that divide 541 them. Future research would determine whether UBI outperforms MBI.

Explanans	Explanandum	Scientific Example	Philosophical Work Discussing Example
Computational Explanations			
Difference of Gaussians	Stereoscopic Vision	Marr (1982); Rodieck (1965)	Bechtel and Shagrir (2015); Egan (2017); Kaplan (2011*); Kaplan and Craver (2011); Rusanen and Lappi (2016); Shagrir (2010); Shapiro (2019)
Exhaustive Search	Recall (Memory)	Sternberg (1969)	Shapiro (2017, 2019)
Gain Field Encoding	Hand-Eye Coordination	Pouget, Deneve, and Duhamel (2002); Pouget and Sejnowski (1997); Shadmehr and Wise (2005); Zipser and Andersen (1988)	Egan (2017); Kaplan (2011*); Rusanen and Lappi (2016); Serban (2015); Shagrir (2006*)

Geon Composition	Object Recognition	Hummel and Biederman (1992)	Buckner (2015*); Povich (2015*); Weiskopf (2011)
Hybrid Computation	Efficiency of Brain	Sarpeshkar (1998)	Chirimuuta (2018)
Inhibitory Feedback	Normalization	Carandini and Heeger (2012)	Chirimuuta (2014); Serban (2015)
Internal Integration	Eye Movement	Seung, Lee, Reis, and Tank (2000)	Egan (2017)
Line Attractor of Choice Axis, Stimuli's Selection Vector	Context- Dependent Decision Making	Mante, Sussillo, Shenoy, and Newsome (2013)	Chirimuuta (2018)
Mapping Non- Coplanar Points to Unique Rigid Configuration	Three- Dimensional Visual Structure of Moving Objects	Ullman (1979)	Egan (2017); Shagrir and Bechtel (2014*)
Optimization of Spatial and Spectral Information Recovery (Gabor Function)	V1 Receptive Fields	Daugman (1985)	Chirimuuta (2014, 2018)
Similarity of Stimulus to Stored Exemplars	Categorization	Kruschke (2008); Love, Medin, and Gureckis (2004)	Buckner (2015*); Povich (2015*); Weiskopf (2011)
Topological Explanations			
Closeness Centrality	Speech and Tonal Processing	Mišić et al. (2018)	Kostić (2020)
Mean Connectivity	Ictogenicity	Helling, Petkov, and Kalitzin (2019)	Kostić and Khalifa (2021)

Motif Frequency	Functional Connectivity	Adachi et al. (2011)	Kostić and Khalifa (2021, manuscript)
Navigation Efficiency, Diffusion Efficiency	Efficiency of Neuronal Communication	Seguin, Razi, and Zalesky (2019)	Kostić (2020)
Network Communicability	Cognitive Control	Gu et al. (2015)	Kostić (2020)
Small-Worldness	Information Propagation	Watts and Strogatz (1998)	Kostić and Khalifa (manuscript)
	Dynam	ical Explanations	
Coupling of Eye and Bodily Movements	Onset of Motor Control	Kelso et al. (1998); Shenoy, Sahani, and Churchland (2013)	Chemero and Silberstein (2008); Favela (2020b); Vernazzani (2019*)
Coupling Ratio	Bimanual Coordination (Relative Phase)	Haken et al. (1985)	Chemero (2000, 2001); Golonka and Wilson (2019*); Kaplan and Craver (2011*); Lamb and Chemero (2014); Meyer (2018); Stepp, Chemero, and Turvey (2011); Zednik (2011*)
Strength of Memory Trace, Salience of Target, Waiting Time, Stance	Infant Reaching (A-not-B Error)	Thelen, Schöner, Scheier, and Smith (2001)	Gervais (2015); Meyer (2018); Povich (forthcoming*); van Eck (2018*); Venturelli (2016);

			Verdejo (2015); Zednik (2011*)
Potassium and Sodium Ion Flows	Neural Excitability	FitzHugh (1961); Hodgkin and Huxley (1952); Nagumo, Arimoto, and Yoshizawa (1962)	Craver and Kaplan (2011*); Favela (2020a, 2020b); Hochstein (2017*); Kaplan and Bechtel (2011*); Kaplan and Craver (2011*); Ross (2015)

542

Table 2. Putatively non-mechanistic explanations discussed by philosophers.

The *explanans* (first column) is the factor that explains. The *explanandum* (second column) is the
 phenomenon to be explained. An asterisk indicates that the author takes the explanation to be
 mechanistic.

546 4.3.1 Assimilation Strategy

547 According to mechanists' assimilation strategy, many putatively non-mechanistic explanations are in 548 fact elliptical mechanistic explanations or "mechanism sketches" (Miłkowski, 2013; Piccinini, 2015; 549 Piccinini & Craver, 2011; Povich, 2015, forthcoming; Zednik, 2011). Thus, when deploying the assimilation strategy, mechanists take computational, topological, and dynamical models to fall short 550 of a (complete) mechanistic explanation, but to nevertheless provide important information about 551 552 such mechanistic explanations. Mechanists have proposed two ways that putatively non-mechanistic 553 explanations can provide mechanistic information, and thereby serve as mechanism sketches. First, 554 putatively non-mechanistic explanations can be *heuristics* for discovering mechanistic explanations. 555 Second, putatively non-mechanistic explanations can *constrain* the space of acceptable mechanistic 556 explanations.

557 An alternative interpretation is possible. The fact that non-mechanistic models assist in the 558 identification of mechanistic explanations does not entail that the former is a species of the latter. 559 Consequently, putatively non-mechanistic explanations can play these two roles with respect to 560 mechanistic explanations without being mere mechanism sketches. In other words, "genuinely" *non-*561 *mechanistic* explanations can guide or constrain the discovery of *mechanistic* explanations. Earlier 562 explanatory pluralists (McCauley, 1986, 1996) already anticipated precursors to this idea, but did not 563 tie it explicitly as a response to mechanists' assimilation strategy.

Moreover, this fits comfortably with our account of scientific explanatory evaluation (SEEing) and hence with UBI. Heuristics of discovery are naturally seen as advancing SEEing's first stage of considering plausible potential explanations. Similarly, since the goal of SEEing is to identify correct explanations and their relationships, it is a consequence of UBI that different kinds of explanations of the related phenomena constrain each other. For instance, suppose that we have two computational explanations of the same phenomenon, and that the key difference between them is that only the first of these is probable given the best mechanistic explanations of that phenomenon.

- 571 Then that counts as a reason to treat the first computational explanation as better than the second.
- 572 Hence, SEEing entails mechanistic explanations can constrain computational explanations.

573 More generally, UBI can capture the same key inter-explanatory relationships that mechanists

- 574 prize without assimilating putatively non-mechanistic explanations to mechanistic explanation.
- 575 Indeed, like many mechanists, UBI suggests that not only do putatively non-mechanistic explanations
- 576 guide and constrain the discovery of mechanistic explanations, but that the converse is also true. 577 (2 + i) + (2 + 2) + (1 + i) + (1 + i)
- 577 (Section 4.3.2.2 provides an example of this.) Parity of reasoning entails that mechanistic 578 explanations should thereby be relegated to mere "computational, topological, and dynamical
- site explanations should thereby be relegated to mere "computational, topological, and dynamical sketches" in these cases, but mechanists must resist this conclusion on pain of contradiction. Since
- 580 UBI captures these important inter-explanatory relationships without broaching the more
- 581 controversial question of assimilation, it need not determine which models are mere sketches of
- 582 adequate explanations. Future research would evaluate whether this is a virtue or a vice.

583 4.3.2 Negative Strategy

- 584 Mechanists' assimilation strategy becomes more plausible than the UBI-inspired alternative if there
- are good grounds for thinking that the criteria that pluralists use to establish putatively non-
- 586 mechanistic explanations as genuine explanations are insufficient. This is the crux of the mechanists'
- 587 negative strategy. As with the assimilation strategy, we suggest that UBI provides a suggestive foil to
- 588 the negative strategy.
- 589 The negative strategy's key move is to identify a set of non-explanatory models that 590 pluralists' criteria would wrongly label as explanatory. Two kinds of non-explanatory models-how-591 possibly and phenomenological models—exemplify this mechanist argument. How-possibly models describe factors that *could* but do not *actually* produce the phenomenon to be explained. For instance, 592 593 most explanations begin as conjectures or untested hypotheses. Those that turn out to be false will be 594 how-possibly explanations. Phenomenological models, which accurately describe or predict the target phenomenon without explaining it, provide a second basis for the negative strategy. 595 596 Paradigmatically, phenomenological models correctly represent non-explanatory correlations 597 between two or more variables. Mechanists claim that pluralist criteria of explanation will wrongly classify some how-possibly and some phenomenological models as correct explanations. By contrast, 598 599 since models that accurately represent mechanisms are "how-actually models," i.e., models that cite 600 explanatory factors responsible for the phenomenon of interest, MBI appears well-positioned to distinguish correct explanations from how-possibly and phenomenological models. 601
- 602 However, UBI can distinguish correct explanations from how-possibly and phenomenological 603 models. Moreover, it can do so in two distinct ways that do not appeal to mechanisms. First, it can do 604 so on what we call *structural* grounds, i.e., by identifying non-mechanistic criteria of explanation that 605 are sufficient for funding the distinction. It can also defuse the negative strategy on what we call 606 *procedural* grounds, i.e., by showing that the procedures and methods that promote understanding 607 also distinguish correct explanations from these non-explanatory models.
- 608 4.3.2.1.Structural Defenses
- 609 We suggest that the following provides a structural defense against the negative strategy:
- 610 If *X* correctly explains *Y*, then the following are true:
- 611 (1) Accuracy Condition: X is an accurate representation, and

612 (2) *Counterfactual Condition*: Had the objects, processes, etc. represented by *X* been different,
 613 then *Y* would have been different.

614 These are only necessary conditions for correct explanations. They are also sufficient for 615 distinguishing correct explanations from how-possibly and phenomenological models but are likely 616 insufficient for distinguishing correct explanations from every other kind of non-explanatory model. 617 Identifying these other models is a useful avenue for future iterations of the negative strategy and 618 responses thereto.

619 Situating this within UBI, these conditions are naturally seen as elaborating the Explanatory 620 Floor, which claims that understanding a phenomenon requires possession of a correct explanation. Crucially, mechanists and pluralists alike widely accept these as requirements on correct 621 622 explanations, though we discuss some exceptions below. Reasons for their widespread acceptance becomes clear with a simple example. Consider a case in which it is hypothesized that taking a 623 624 certain medication (X) explains recovery from an illness (Y). If it were discovered that patients had 625 not taken the medication, then this hypothesis would violate the accuracy condition. Intuitively, it 626 would not be a correct explanation, but it would be a how-possibly model.

More generally, how-possibly models are correct explanations *modulo* satisfaction of the accuracy condition. Consequently, pluralists can easily preserve this distinction without appealing to mechanisms; accuracy is sufficient. Just as mechanisms can be either accurately or inaccurately represented, so too can computations, topological structures, and system dynamics be either accurately or inaccurately represented. Similarly, just as inaccurate mechanistic models can be howpossibly models but cannot be correct explanations, so too can inaccurate computational, topological, and dynamical models be how-possibly models but cannot be how-actually models.

Analogously, the counterfactual condition preserves the distinction between correct explanations 634 635 and phenomenological models. Suppose that our hypothesis about recovery is confounded by the fact 636 that patients' recovery occurred two weeks after the first symptoms, and that this is the typical 637 recovery time for anyone with the illness in question, regardless of whether they take medication. Barring extenuating circumstances, e.g., that the patients are immunocompromised, these facts would 638 639 seem to cast doubt upon the claim that the medication made a difference to their recovery. In other 640 words, they cast doubt on the following counterfactual: had a patient not taken the medication, then 641 that patient would not have recovered when she did. Consequently, the hypothesis about the 642 medication explaining recovery violates the counterfactual condition. Moreover, the hypothesis does 643 not appear to be correct, but would nevertheless describe the patients' situation, i.e., it would be a 644 phenomenological model.

More generally, phenomenological models are correct explanations *modulo* satisfaction of the counterfactual condition. Just as a mechanistic model may accurately identify interacting parts of a system that correlate with but do not explain its behavior, a non-mechanistic model may accurately identify computational processes, topological structures, and dynamical properties of a system that correlate with but do not explain its behavior. In both cases, the counterfactual condition accounts for the models' explanatory shortcomings; no appeal to mechanisms is needed.

651 4.3.2.2.Procedural Defenses

- 652 Admittedly, structural defenses against the negative strategy are not unique to UBI; other pluralists
- 653 who are agnostic about UBI have invoked them in different ways. By contrast, our second
- 654 procedural defense against the negative strategy is part and parcel to UBI. Procedural defenses show

- 655 that the procedures that promote understanding also distinguish correct explanations from how-
- 656 possibly and phenomenological models.

The Scientific Knowledge Principle characterizes the key procedures that simultaneously 657 658 promote understanding and distinguish correct explanations from these non-explanatory models. 659 Recall that SEEing consists of three stages: considering plausible potential explanations of a phenomenon, *comparing* them using the best available methods, and forming *commitments* to 660 explanatory models based on these comparisons. This provides a procedural defense against the 661 662 negative strategy. How-possibly and phenomenological models will only be acceptable in the first stage of SEEing: prior to their deficiencies being discovered, they frequently deserve consideration 663 664 as possible explanations of a phenomenon. By contrast, correct explanations must "survive" the 665 remaining stages of SEEing: they must pass certain empirical tests at the stage of comparison such that they are acceptable at the stage of commitment. Indeed, it is often through SEEing that scientists 666 come to distinguish correct explanations from how-possibly and phenomenological models. 667

668 Crucially, consideration is most effective when it does not prejudge what makes something genuinely explanatory. This minimizes the possibility of missing out on a fruitful hypothesis. 669 670 Consequently, both mechanistic and non-mechanistic explanations should be included at this initial 671 stage of SEEing. However, our procedural defense supports pluralism only if some computational, 672 topological, or dynamical explanations are acceptable in light of rigorous explanatory comparisons. 673 As we see it, this is a strength of our procedural defense, for it uses the empirical resources of our 674 best science to adjudicate debates between mechanists and pluralists that often appear intractable 675 from the philosophical armchair.

Nevertheless, we can point to an important kind of explanatory comparison-which we call 676 677 control-and-contrast-that deserves greater philosophical and scientific attention when considering explanatory integration in the cognitive sciences. Control-and-contrast proceeds as follows. Let X_1 678 679 and X_2 be two potential explanations of Y under consideration. Next, run two controlled experiments: 680 one in which the explanatory factors in X_1 are absent but those in X_2 are present and the second in 681 which the explanatory factors in X_1 are present but those in X_2 are absent. If Y is only present in the 682 first experiment, then the pair of experiments suggests that X_2 is a better explanation of Y than X_1 . 683 Conversely, if Y is only present in the second experiment, the pair of experiments suggests that X_1 is 684 a better explanation of Y than X_2 . If Y is present in both experiments, the experiments are 685 inconclusive. If Y is absent in both experiments, then the experiments suggest that the combination of 686 X_1 and X_2 better explains Y than either X_1 or X_2 does in isolation. Since we suggest that both 687 mechanistic and non-mechanistic explanations should be considered and thereby play the roles of X_1 688 and X_2 , we also suggest that which of these different kinds of explanations is correct for a given 689 phenomenon Y should frequently be determined by control-and-contrast.

690 In some cases, scientists are only interested in controlling-and-contrasting explanations of the 691 same kind. However, even in these cases, the controls are often best described in terms of other kinds 692 of explanation. For instance, as discussed above, the default mode network mechanistically explains 693 mental simulations involved in episodic memory. By contrast, when mental simulations involve 694 semantic memory, inferior temporal and lateral occipital regions play a more pronounced role (Parikh 695 et al., 2018). Both episodic and semantic memory are functional or computational concepts that can 696 figures as controls in different experiments designed to discover which of these mechanisms explains 697 a particular kind of mental simulation. Less common is controlling-and-contrasting explanations of 698 different kinds. Perhaps this is a lacuna in current research. Alternatively, it may turn out that

699 different kinds of explanation rarely compete and are more amenable to integration in the ways700 outlined above.

701 The procedural defense complements the structural defense in two ways. First, not all 702 pluralists accept the accuracy condition. Their motivations for this are twofold. First, given that 703 science is a fallible enterprise, our best explanations today are likely to be refuted. Second, many 704 explanations invoke idealizations, i.e., known inaccuracies that nevertheless enhance understanding. 705 The procedural defense does not require the accuracy condition but can still preserve the distinction 706 between correct explanations and non-explanatory models. Instead, the procedural defense only 707 requires that correct explanations be acceptable on the basis of the best available scientific methods 708 and evidence.

- 709 Second, tests such as control-and-contrast regiment the subjunctive conditionals that 710 characterize the counterfactual condition. In evaluating counterfactuals, it is notoriously difficult to 711 identify what must be held constant, what can freely vary without altering the truth-value of the 712 conditional, and what must vary in order to determine the truth-value of the conditional. Our account 713 of explanatory evaluation points to important constraints on this process. Suppose that we are 714 considering two potential explanations X_i and X_j of some phenomenon Y. To compare these models, 715 we will be especially interested in counterfactuals such as, "Had the value of X_i been different (but 716 the value of X_i had remained the same), then the value of Y would have been different," and also, 717 "Had the value of X_i been different (but the value of X_i had remained the same), then the value of Y
- would have been the same." These are precisely the kinds of counterfactuals that will be empirically
- 719 supported or refuted by control-and-contrast.

720 **5** Conclusion

721 Fruitful connections between the philosophy and science of understanding can be forged. In a

- naturalized epistemology of understanding, philosophical claims about various forms of explanatory
- and counterfactual reasoning are empirically constrained by scientific tests and explanations. By
- contrast, in understanding-based integration, the philosophy of understanding contributes to the
- science of understanding by providing broad methodological prescriptions as to how diverse
- explanations can be woven together. Specifically, understanding-based integration includes
- identification of inter-explanatory relationships, consideration of different kinds of explanations, and
- evaluation of these explanations using methods such as control-and-contrast. As our suggestions have
 been of a preliminary character, we hope that future collaborations between philosophers and
- 129 Deen of a premimary character, we nope that future collaborations between philosophe
 720 scientists will advance our understanding of understanding
- radia scientists will advance our understanding of understanding.
- Adachi, Y., Osada, T., Sporns, O., Watanabe, T., Matsui, T., Miyamoto, K., & Miyashita, Y. (2011).
 Functional Connectivity between Anatomically Unconnected Areas Is Shaped by Collective
 Network-Level Effects in the Macaque Cortex. *Cerebral Cortex*, 22(7), 1586-1592.
 doi:10.1093/cercor/bhr234
- Baumberger, C. (2014). Types of understanding: Their nature and their relation to knowledge.
 Conceptus, 40(98), 67–88.
- Baumberger, C. (2019). Explicating Objectual Understanding: Taking Degrees Seriously. *Journal for General Philosophy of Science*, 50(3), 367-388. doi:10.1007/s10838-019-09474-6
- Baumberger, C., Beisbart, C., & Brun, G. (2016). What is understanding? An overview of recent
 debates in epistemology and philosophy of science. In S. R. Grimm, C. Baumberger, & S.

- Ammon (Eds.), *Explaining understanding: New perspectives from epistemology and philosophy of science* (pp. 1-34). New York: Routledge.
- Baumberger, C., & Brun, G. (2017). Dimensions of Objectual Understanding. In S. G. Christoph
 Baumberger & S. Ammon (Eds.), *Explaining understanding: New perspectives from epistemology and philosophy of science* (pp. 165-189): Routledge.
- Bechtel, W., & Richardson, R. C. (1993). *Discovering complexity : decomposition and localization as strategies in scientific research*. Princeton, N.J.: Princeton University Press.
- Bechtel, W., & Shagrir, O. (2015). The Non-Redundant Contributions of Marr's Three Levels of
 Analysis for Explaining Information-Processing Mechanisms. *Topics in Cognitive Science*,
 750 7(2), 312-322. doi:<u>https://doi.org/10.1111/tops.12141</u>
- Bokulich, A. (2011). How scientific models can explain. *Synthese*, *180*(1), 33-45.
 doi:10.1007/s11229-009-9565-1
- 753 Buckner, C. (2015). Functional kinds: a skeptical look. *Synthese*, 192(12), 3915-3942.
- Bullmore, E., & Sporns, O. (2012). The economy of brain network organization. *Nature Reviews Neuroscience*, 13(5), 336-349. doi:10.1038/nrn3214
- Carandini, M., & Heeger, D. J. (2012). Normalization as a canonical neural computation. *Nature Reviews Neuroscience*, 13(1), 51-62. doi:10.1038/nrn3136
- Carter, J. A., & Gordon, E. C. (2014). Objectual understanding and the value problem. *American philosophical quarterly*, *51*(1), 1-13.
- Chemero, A. (2000). Anti-Representationalism and the Dynamical Stance. *Philosophy of Science*,
 67(4), 625-647. doi:10.1086/392858
- Chemero, A. (2001). Dynamical explanation and mental representations. *Trends in cognitive sciences*, 5(4), 141-142.
- 764 Chemero, A. (2009). Radical embodied cognitive science. Cambridge, Mass.: MIT Press.
- Chemero, A., & Silberstein, M. (2008). After the Philosophy of Mind: Replacing Scholasticism with
 Science. *Philosophy of Science*, 75(1), 1-27. doi:10.1086/587820
- Cheng, P. W., & Holyoak, K. J. (1985). Pragmatic reasoning schemas. *Cognitive psychology*, *17*(4), 391-416. doi:<u>https://doi.org/10.1016/0010-0285(85)90014-3</u>
- Cheng, P. W., Holyoak, K. J., Nisbett, R. E., & Oliver, L. M. (1986). Pragmatic versus syntactic approaches to training deductive reasoning. *Cognitive psychology*, *18*(3), 293-328.
 doi:<u>https://doi.org/10.1016/0010-0285(86)90002-2</u>
- Chirimuuta, M. (2014). Minimal models and canonical neural computations: the distinctness of
 computational explanation in neuroscience. *Synthese*, *191*(2), 127-153. doi:10.1007/s11229013-0369-y
- Chirimuuta, M. (2018). Explanation in Computational Neuroscience: Causal and Non-causal. *The British Journal for the Philosophy of Science*, 69(3), 849-880. doi:10.1093/bjps/axw034
- 777 Chollet, F. (2019). On the measure of intelligence. *arXiv preprint arXiv:1911.01547*.
- 778 Craver, C. F. (2007). *Explaining the brain: mechanisms and the mosaic unity of neuroscience*.
 779 Oxford: Clarendon Press.

- Craver, C. F. (2014). The Ontic Account of Scientific Explanation. In I. M. Kaiser, R. O. Scholz, D.
 Plenge, & A. Hüttemann (Eds.), *Explanation in the Special Sciences: The Case of Biology and History* (pp. 27-52). Dordrecht: Springer Netherlands.
- 783 Craver, C. F., & Kaplan, D. M. (2011). Towards a Mechanistic Philosophy of Neuroscience. In S.
 784 French & J. Saatsi (Eds.), *Continuum Companion to the Philosophy of Science* (pp. 268):
 785 Continuum.
- Craver, C. F., & Kaplan, D. M. (2020). Are More Details Better? On the Norms of Completeness for
 Mechanistic Explanations. *The British Journal for the Philosophy of Science*, *71*(1), 287-319.
 doi:10.1093/bjps/axy015
- Craver, C. F., & Tabery, J. (2019). Mechanisms in Science. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy (Summer 2019 Edition)*. URL =
 https://plato.stanford.edu/archives/sum2019/entries/science-mechanisms/>.
- 792 Cummins, R. C. (1975). Functional analysis. *Journal of philosophy*, 72(20), 741-765.
- 793 Cummins, R. C. (1983). The nature of psychological explanation. Cambridge, Mass.: MIT Press.
- Cummins, R. C. (2000). How does it work?" versus" what are the laws?": Two conceptions of
 psychological explanation. *Explanation and cognition*, 117-144.
- Darrason, M. (2018). Mechanistic and topological explanations in medicine: the case of medical
 genetics and network medicine. *Synthese*, 195(1), 147-173. doi:10.1007/s11229-015-0983-y
- Daugman, J. G. (1985). Uncertainty relation for resolution in space, spatial frequency, and orientation
 optimized by two-dimensional visual cortical filters. *Journal of the Optical Society of America A*, 2(7), 1160-1169. doi:10.1364/JOSAA.2.001160
- Boli Davies, P. C. (1990). Why is the physical world so comprehensible. In W. Zurek (Ed.), *Complexity, entropy and the physics of information* (pp. 61-70): Addison-Wesley Publishing Company.
- 803 De Regt, H. W. (2017). Understanding scientific understanding. New York: Oxford University Press.
- Bellsén, F. (2020). Beyond Explanation: Understanding as Dependency Modelling. *The British Journal for the Philosophy of Science*, 71(4), 1261-1286. doi:10.1093/bjps/axy058
- Bowhurst, J. (2018). Individuation without Representation. *The British Journal for the Philosophy of Science*, 69(1), 103-116. doi:10.1093/bjps/axw018
- Egan, F. (2017). Function-Theoretic Explanation and the Search for Neural Mechanisms. In D. M.
 Kaplan (Ed.), *Explanation and Integration in Mind and Brain Science* (pp. 145-163). Oxford:
 Oxford University Press.
- 811 Elgin, C. Z. (2004). True enough. *Philosophical issues*, *14*(1), 113-131. Retrieved from
 812 <u>http://dx.doi.org/10.1111/j.1533-6077.2004.00023.x</u>
- 813 Elgin, C. Z. (2017). *True enough*. Cambridge, MA: MIT Press.
- Evans, J. S. B. T. (2011). Dual-process theories of reasoning: Contemporary issues and
 developmental applications. *Developmental Review*, *31*(2), 86-102.
 doi:<u>https://doi.org/10.1016/j.dr.2011.07.007</u>
- Evans, J. S. B. T. (2012). Dual-process theories of deductive reasoning: Facts and fallacies. In *The Oxford handbook of thinking and reasoning*. (pp. 115-133). New York, NY, US: Oxford
 University Press.

- Favela, L. H. (2020a). The dynamical renaissance in neuroscience. *Synthese*. doi:10.1007/s11229020-02874-y
- Favela, L. H. (2020b). Dynamical systems theory in cognitive science and neuroscience. *Philosophy compass*, 15(8), e12695. doi:<u>https://doi.org/10.1111/phc3.12695</u>
- FitzHugh, R. (1961). Impulses and Physiological States in Theoretical Models of Nerve Membrane.
 Biophysical Journal, 1(6), 445-466. doi:<u>https://doi.org/10.1016/S0006-3495(61)86902-6</u>
- Fodor, J. A. (1968). *Psychological Explanation: An Introduction To The Philosophy Of Psychology*.
 New York: Random House.
- Fresco, N., & Miłkowski, M. (2021). Mechanistic Computational Individuation without Biting the
 Bullet. *The British Journal for the Philosophy of Science*, 72(2), 431-438.
 doi:10.1093/bjps/axz005
- 831 Friedman, M. (1974). Explanation and scientific understanding. *Journal of philosophy*, 71(1), 5-19.
- Gervais, R. (2015). Mechanistic and non-mechanistic varieties of dynamical models in cognitive
 science: explanatory power, understanding, and the 'mere description' worry. *Synthese*, *192*(1), 43-66. doi:10.1007/s11229-014-0548-5
- Glennan, S. (2017). *The new mechanical philosophy* (First edition. ed.). Oxford: Oxford University
 Press.
- Goel, V., Buchel, C., Frith, C., & Dolan, R. J. (2000). Dissociation of Mechanisms Underlying
 Syllogistic Reasoning. *NeuroImage*, *12*(5), 504-514.
 doi:https://doi.org/10.1006/nimg.2000.0636
- Golonka, S., & Wilson, A. D. (2019). Ecological mechanisms in cognitive science. *Theory & Psychology*, 29(5), 676-696. doi:10.1177/0959354319877686
- 842 Gopnik, A. (1998). Explanation as Orgasm. *Minds and Machines*, 8(1), 101-118.
 843 doi:10.1023/A:1008290415597
- Gordon, E. C. (2017). Understanding in Epistemology. In *Internet Encyclopedia of Philosophy*.
 <u>https://iep.utm.edu/understa/</u>.
- 846 Greco, J. (2013). Episteme: knowledge and understanding. In K. Timpe & C. A. Boyd (Eds.), *Virtues* 847 *and their vices* (pp. 285-301). Oxford: Oxford University Press.
- Grimm, S. R. (2010). The goal of understanding. *Studies in the history and philosophy of science*,
 41(4), 337-344. doi:10.1016/j.shpsa.2010.10.006
- Grimm, S. R. (2014). Understanding as knowledge of causes. In A. Fairweather (Ed.), *Virtue epistemology naturalized* (Vol. 366, pp. 329-345). Dordecht: Springer International
 Publishing.
- 853 Grimm, S. R. (2021). Understanding. In E. N. Zalta (Ed.), *Stanford encyclopedia of philosophy*.
- Grünwald, P. (2004). A tutorial introduction to the minimum description length principle. *arXiv preprint math/0406077*.
- Gu, S., Pasqualetti, F., Cieslak, M., Telesford, Q. K., Yu, A. B., Kahn, A. E., ... Bassett, D. S.
 (2015). Controllability of structural brain networks. *Nature Communications*, 6(1), 8414.
 doi:10.1038/ncomms9414
- Haken, H., Kelso, J. A. S., & Bunz, H. (1985). A theoretical model of phase transitions in human
 hand movements. *Biological Cybernetics*, *51*(5), 347-356. doi:10.1007/BF00336922

- Hannon, M. (2021). RECENT WORK IN THE EPISTEMOLOGY OF UNDERSTANDING.
 American philosophical quarterly, 58(3), 269-290. doi:10.2307/48616060
- Helling, R. M., Petkov, G. H., & Kalitzin, S. N. (2019). *Expert system for pharmacological epilepsy treatment prognosis and optimal medication dose prescription: computational model and clinical application*. Paper presented at the Proceedings of the 2nd International Conference
 on Applications of Intelligent Systems, <u>https://doi.org/10.1145/3309772.3309775</u>.
 <u>https://doi.org/10.1145/3309772.3309775</u>
- 868 Hills, A. (2015). Understanding why. *Noûs*, 49(2), 661-688. doi:10.1111/nous.12092
- Hitchcock, C. R., & Woodward, J. (2003). Explanatory generalizations, part II: plumbing explanatory
 depth. *Noûs*, *37*(2), 181-199.
- Hochstein, E. (2016). One mechanism, many models: a distributed theory of mechanistic
 explanation. *Synthese*, *193*(5), 1387-1407.
- Hochstein, E. (2017). Why one model is never enough: a defense of explanatory holism. *Biology & Philosophy*, *32*(6), 1105-1125. doi:10.1007/s10539-017-9595-x
- Hodgkin, A. L., & Huxley, A. F. (1952). A quantitative description of membrane current and its
 application to conduction and excitation in nerve. *The Journal of Physiology*, *117*(4), 500544. doi:https://doi.org/10.1113/jphysiol.1952.sp004764
- Holyoak, K. J., & Cheng, P. W. (1995). Pragmatic reasoning with a point of view. *Thinking & Reasoning*, 1(4), 289-313. doi:10.1080/13546789508251504
- Hopkins, E. J., Weisberg, D. S., & Taylor, J. C. V. (2016). The seductive allure is a reductive allure:
 People prefer scientific explanations that contain logically irrelevant reductive information.
 Cognition, 155, 67-76. doi:https://doi.org/10.1016/j.cognition.2016.06.011
- Hummel, J. E., & Biederman, I. (1992). Dynamic binding in a neural network for shape recognition. *Psychological Review*, 99(3), 480-517. doi:10.1037/0033-295X.99.3.480
- Humphreys, P. (1993). Greater Unification Equals Greater Understanding? *Analysis*, *53*(3), 183-188.
 doi:10.2307/3328470
- Huneman, P. (2018). Outlines of a theory of structural explanations. *Philosophical studies*, 175(3),
 665-702. doi:10.1007/s11098-017-0887-4
- Illari, P. M., & Williamson, J. (2010). Function and organization: comparing the mechanisms of
 protein synthesis and natural selection. *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological and Biomedical Sciences*, 41(3), 279-291.
 doi:https://doi.org/10.1016/j.shpsc.2010.07.001
- Janssen, A., Klein, C., & Slors, M. (2017). What is a cognitive ontology, anyway? *Philosophical Explorations*, 20(2), 123-128. doi:10.1080/13869795.2017.1312496
- Johnson-Laird, P. N. (1995). Mental models, deductive reasoning, and the brain. In M. S. Gazzaniga
 (Ed.), *The Cognitive Neurosciences* (pp. 999--1008): MIT Press.
- Johnson-Laird, P. N. (2010). Mental models and human reasoning. *Proceedings of the National Academy of Sciences*, 107(43), 18243-18250. doi:10.1073/pnas.1012933107
- Kaplan, D. M. (2011). Explanation and description in computational neuroscience. *Synthese*, 183(3),
 339-373.

- Kaplan, D. M. (2017). *Explanation and integration in mind and brain science* (First edition. ed.).
 Oxford, United Kingdom: Oxford University Press.
- Kaplan, D. M., & Bechtel, W. (2011). Dynamical Models: An Alternative or Complement to
 Mechanistic Explanations? *Topics in Cognitive Science*, 3(2), 438-444.
 doi:<u>https://doi.org/10.1111/j.1756-8765.2011.01147.x</u>
- Kaplan, D. M., & Craver, C. F. (2011). The Explanatory Force of Dynamical and Mathematical
 Models in Neuroscience: A Mechanistic Perspective. *Philosophy of Science*, 78(4), 601-627.
 doi:10.1086/661755
- Kelp, C. (2015). Understanding phenomena. Synthese, 192(12), 3799–3816. doi:10.1007/s11229-014-0616-x
- Kelso, J. A. S., Fuchs, A., Lancaster, R., Holroyd, T., Cheyne, D., & Weinberg, H. (1998). Dynamic
 cortical activity in the human brain reveals motor equivalence. *Nature*, *392*(6678), 814-818.
 doi:10.1038/33922
- Keren, G., & Schul, Y. (2009). Two Is Not Always Better Than One:A Critical Evaluation of TwoSystem Theories. *Perspectives on Psychological Science*, 4(6), 533-550. doi:10.1111/j.17456924.2009.01164.x
- Khalifa, K. (2012). Inaugurating understanding or repackaging explanation? *Philosophy of Science*,
 79(1), 15-37. Retrieved from http://www.jstor.org/stable/10.1086/663235
- Synthese, 190(6), 1153-1171.
 doi:10.1007/s11229-011-9886-8
- Khalifa, K. (2013b). The role of explanation in understanding. *British Journal for the Philosophy of Science*, 64(1), 161-187. doi:10.1093/bjps/axr057
- Khalifa, K. (2017). Understanding, Explanation, and Scientific Knowledge. Cambridge: Cambridge
 University Press.
- Khalifa, K. (2019). Is *Verstehen* Scientific Understanding? *Philosophy of the Social Sciences*, 49(4),
 282-306. doi:10.1177/0048393119847104
- Khalifa, K. (forthcoming). Should Friends and Frenemies of Understanding be Friends? Discussing
 de Regt. In K. Khalifa, I. Lawler, & E. Shech (Eds.), *Scientific Understanding and Representation: Modeling in the Physical Sciences*. London: Routledge.
- Kitcher, P. (1989). Explanatory unification and the causal structure of the world. In P. Kitcher & W.
 C. Salmon (Eds.), *Scientific explanation* (Vol. XIII, pp. 410-506). Minneapolis: University of Minnesota Press.
- Kohlberg, L. (1958). *The development of modes of moral thinking and choice in the years 10 to 16.*(PhD). University of Chicago, Chicago. /z-wcorg/ database.
- Kon, E., & Lombrozo, T. (2019). Scientific Discovery and the Human Drive to Explain. In D. A.
 Wilkenfeld & R. Samuels (Eds.), *Advances in Experimental Philosophy of Science* (pp. 15).
 London: Routledge.
- Korb, K. B. (2004). Introduction: Machine Learning as Philosophy of Science. *Minds and Machines*, *14*(4), 433-440. doi:10.1023/B:MIND.0000045986.90956.7f
- Koslowski, B., Marasia, J., Chelenza, M., & Dublin, R. (2008). Information becomes evidence when
 an explanation can incorporate it into a causal framework. *Cognitive Development*, 23(4),
 472-487. doi:<u>https://doi.org/10.1016/j.cogdev.2008.09.007</u>

- Wostić, D. (2018). The topological realization. Synthese, 195(1), 79-98. doi:10.1007/s11229-016 1248-0
- Kostić, D. (2020). General theory of topological explanations and explanatory asymmetry.
 Philosophical Transactions of the Royal Society B: Biological Sciences, 375(1796),
 20190321. doi:doi:10.1098/rstb.2019.0321
- Kostić, D., & Khalifa, K. (2021). The directionality of topological explanations. *Synthese*.
 doi:10.1007/s11229-021-03414-y
- 950 Kostić, D., & Khalifa, K. (manuscript). Decoupling Topological Explanation from Mechanisms.
- Kroger, J. K., Nystrom, L. E., Cohen, J. D., & Johnson-Laird, P. N. (2008). Distinct neural substrates
 for deductive and mathematical processing. *Brain Research*, *1243*, 86-103.
 doi:<u>https://doi.org/10.1016/j.brainres.2008.07.128</u>
- Kruschke, J. K. (2008). Models of categorization. In R. Sun (Ed.), *The Cambridge Handbook of Computational Psychology* (pp. 267--301): Cambridge University Press.
- Kuorikoski, J., & Ylikoski, P. (2015). External representations and scientific understanding.
 Synthese, 192(12), 3817-3837. doi:10.1007/s11229-014-0591-2
- Kvanvig, J. L. (2003). *The value of knowledge and the pursuit of understanding*. Cambridge:
 Cambridge University Press.
- Lamb, M., & Chemero, A. (2014). Structure and application of dynamical models in cognitive
 science. Paper presented at the Proceedings of the annual meeting of the cognitive science
 society.
- Lange, M. (2017). Because without cause : non-causal explanation in science and mathematics. New
 York: Oxford University Press.
- Latora, V., & Marchiori, M. (2001). Efficient Behavior of Small-World Networks. *Physical Review Letters*, 87(19), 198701. doi:10.1103/PhysRevLett.87.198701
- Le Bihan, S. (2016). Enlightening Falsehoods: A modal view of scientific understanding. In S. R.
 Grimm, C. Baumberger, & S. Ammon (Eds.), *Explaining understanding: New perspectives* from epistemology and philosophy of science (pp. 111-136). London: Routledge.
- Levy, A. (2014). What was Hodgkin and Huxley's Achievement? *The British Journal for the Philosophy of Science*, 65(3), 469-492. doi:10.1093/bjps/axs043
- Li, M., & Vitányi, P. (2008). An introduction to Kolmogorov complexity and its applications (Vol. 3): Springer.
- Lombrozo, T. (2006). The structure and function of explanations. *Trends in cognitive sciences*,
 10(10), 464-470. doi:<u>https://doi.org/10.1016/j.tics.2006.08.004</u>
- Lombrozo, T., & Wilkenfeld, D. (2019). Mechanistic versus Functional Understanding. In S. R.
 Grimm (Ed.), *Varieties of Understanding* (pp. 209-230): Oxford University Press.
- Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). SUSTAIN: A Network Model of Category
 Learning. *Psychological Review*, 111(2), 309-332.
- Machamer, P., Darden, L., & Craver, C. F. (2000). Thinking about mechanisms. *Philosophy of Science*, 67, 1-25.

- Mante, V., Sussillo, D., Shenoy, K. V., & Newsome, W. T. (2013). Context-dependent computation
 by recurrent dynamics in prefrontal cortex. *Nature*, 503(7474), 78-84.
 doi:10.1038/nature12742
- Marr, D. (1982). Vision : a computational investigation into the human representation and
 processing of visual information. San Francisco: W.H. Freeman.
- McCauley, R. N. (1986). Intertheoretic Relations and the Future of Psychology. *Philosophy of Science*, 53(2), 179-199. Retrieved from http://www.jstor.org/stable/187691
- McCauley, R. N. (1996). Explanatory pluralism and the coevolution of theories in science. In R. N.
 McCauley (Ed.), *The Churchlands and Their Critics* (pp. 17-47): Blackwell Publishers.
- Meyer, R. (2018). The Nonmechanistic Option: Defending Dynamical Explanation. *British Journal for the Philosophy of Science*, 0-0.
- 993 Miłkowski, M. (2013). Explaining the Computational Mind: MIT Press.
- Miłkowski, M., & Hohol, M. (2020). Explanations in cognitive science: unification versus pluralism.
 Synthese. doi:10.1007/s11229-020-02777-y
- Mišić, B., Betzel, R. F., Griffa, A., de Reus, M. A., He, Y., Zuo, X.-N., ... Zatorre, R. J. (2018).
 Network-Based Asymmetry of the Human Auditory System. *Cerebral Cortex*, 28(7), 2655-2664. doi:10.1093/cercor/bhy101
- Nagumo, J., Arimoto, S., & Yoshizawa, S. (1962). An Active Pulse Transmission Line Simulating
 Nerve Axon. *Proceedings of the IRE, 50*(10), 2061-2070. doi:10.1109/JRPROC.1962.288235
- Newman, M. (2012). An inferential model of scientific understanding. *International studies in the philosophy of science*, 26(1), 1-26. doi:10.1080/02698595.2012.653118
- Newman, M. (2013). Refining the inferential model of scientific understanding. *International studies in the philosophy of science*, 27(2), 173-197. doi:10.1080/02698595.2013.813253
- Newman, M. (2015). Theoretical Understanding in Science. *The British Journal for the Philosophy of Science*, 68(2), 571-595. doi:10.1093/bjps/axv041
- Operskalski, J. T., & Barbey, A. K. (2017). Cognitive neuroscience of causal reasoning. In *The Oxford handbook of causal reasoning*. (pp. 217-242). New York, NY, US: Oxford University
 Press.
- 1010 Osman, M. (2004). An evaluation of dual-process theories of reasoning. *Psychonomic Bulletin & Review*, 11(6), 988-1010. doi:10.3758/BF03196730
- 1012 Osman, M. (2014). Reasoning research: Where was it going? Where is it now? Where will it be
 1013 going? In *New approaches in reasoning research*. (pp. 104-123). New York, NY, US:
 1014 Psychology Press.
- Parikh, N., Ruzic, L., Stewart, G. W., Spreng, R. N., & De Brigard, F. (2018). What if? Neural activity underlying semantic and episodic counterfactual thinking. *NeuroImage*, *178*, 332-345. doi:https://doi.org/10.1016/j.neuroimage.2018.05.053
- 1018 Petersen, S. (ms). Explanation as compression.
- Piaget, J. (1952). *The origins of intelligence in children* (M. Cook, Trans.). New York: W. W. Norton
 & Co.
- 1021 Piccinini, G. (2006). Computational explanation in neuroscience. *Synthese*, 153(3), 343-353.

- 1022 Piccinini, G. (2015). *Physical computation: a mechanistic account*. Oxford: Oxford University Press.
- Piccinini, G., & Craver, C. (2011). Integrating psychology and neuroscience: functional analyses as
 mechanism sketches. *Synthese*, 183(3), 283-311. doi:10.1007/s11229-011-9898-4
- Poldrack, R. A., & Yarkoni, T. (2016). From Brain Maps to Cognitive Ontologies: Informatics and the Search for Mental Structure. *Annual Review of Psychology*, 67(1), 587-612.
 doi:10.1146/annurev-psych-122414-033729
- Potochnik, A. (2017). *Idealization and the aims of science*. Chicago: The University of Chicago
 Press.
- Pouget, A., Deneve, S., & Duhamel, J.-R. (2002). A computational perspective on the neural basis of
 multisensory spatial representations. *Nature Reviews Neuroscience*, 3(9), 741-747.
 doi:10.1038/nrn914
- Pouget, A., & Sejnowski, T. J. (1997). Spatial Transformations in the Parietal Cortex Using Basis
 Functions. *Journal of Cognitive Neuroscience*, 9(2), 222-237. doi:10.1162/jocn.1997.9.2.222
- Povich, M. (2015). Mechanisms and Model-Based Functional Magnetic Resonance Imaging.
 Philosophy of Science, 82(5), 1035-1046. doi:10.1086/683438
- Povich, M. (forthcoming). Mechanistic Explanation in Psychology. In H. Stam & H. L. De Jong
 (Eds.), *The SAGE Handbook of Theoretical Psychology*. London: Sage.
- Price, C. J., & Friston, K. J. (2005). Functional ontologies for cognition: The systematic definition of
 structure and function. *Cognitive Neuropsychology*, 22(3-4), 262-275.
 doi:10.1080/02643290442000095
- Pritchard, D. (2009). Safety-based epistemology: whither now? *Journal of philosophical research*,
 34, 33-45.
- Rathkopf, C. (2018). Network representation and complex systems. *Synthese*, *195*(1), 55-78.
 doi:10.1007/s11229-015-0726-0
- Rice, C. (2015). Moving Beyond Causes: Optimality Models and Scientific Explanation. *Noûs*, 49(3),
 589-615. doi:10.1111/nous.12042
- Rodieck, R. W. (1965). Quantitative analysis of cat retinal ganglion cell response to visual stimuli.
 Vision Research, 5(12), 583-601. doi:<u>https://doi.org/10.1016/0042-6989(65)90033-7</u>
- Ross, L. N. (2015). Dynamical Models and Explanation in Neuroscience. *Philosophy of Science*,
 82(1), 32-54.
- Ross, L. N. (2020). Distinguishing topological and causal explanation. *Synthese*.
 doi:10.1007/s11229-020-02685-1
- 1054 Rusanen, A.-M., & Lappi, O. (2016). On computational explanations. *Synthese*, 193(12), 3931-3949.
- Rysiew, P. (2020). Naturalism in Epistemology. In E. N. Zalta (Ed.), *Stanford Encyclopedia of Philosophy*.
- Sarpeshkar, R. (1998). Analog Versus Digital: Extrapolating from Electronics to Neurobiology.
 Neural Computation, 10(7), 1601-1638. doi:10.1162/089976698300017052
- Schank, R. C. (1986). *Explanation patterns : understanding mechanically and creatively*. Hillsdale,
 N.J.: L. Erlbaum Associates.

- Searle, J. R. (1980). Minds, brains, and programs. *Behavioral and Brain Sciences*, 3(3), 417-424.
 doi:10.1017/S0140525X00005756
- Seguin, C., Razi, A., & Zalesky, A. (2019). Inferring neural signalling directionality from undirected
 structural connectomes. *Nature Communications*, 10(1), 4289. doi:10.1038/s41467-019 12201-w
- Serban, M. (2015). The scope and limits of a mechanistic view of computational explanation.
 Synthese, 192(10), 3371-3396.
- Seung, H. S., Lee, D. D., Reis, B. Y., & Tank, D. W. (2000). Stability of the Memory of Eye Position in a Recurrent Network of Conductance-Based Model Neurons. *Neuron*, 26(1), 259-271. doi:<u>https://doi.org/10.1016/S0896-6273(00)81155-1</u>
- Shadmehr, R., & Wise, S. P. (2005). *The computational neurobiology of reaching and pointing : a foundation for motor learning*. Cambridge: MIT Press.
- Shagrir, O. (2006). Why we view the brain as a computer. *Synthese*, *153*(3), 393-416.
 doi:10.1007/s11229-006-9099-8
- Shagrir, O. (2010). Marr on Computational-Level Theories. *Philosophy of Science*, 77(4), 477-500.
 doi:10.1086/656005
- 1077 Shagrir, O., & Bechtel, W. (2014). Marr's Computational Level and Delineating Phenomena.
- Shapiro, L. (2017). Mechanism or Bust? Explanation in Psychology. *The British Journal for the Philosophy of Science*, 68(4), 1037-1059. doi:10.1093/bjps/axv062
- Shapiro, L. (2019). A tale of two explanatory styles in cognitive psychology. *Theory & Psychology*, 29(5), 719-735. doi:10.1177/0959354319866921
- Shenoy, K. V., Sahani, M., & Churchland, M. M. (2013). Cortical Control of Arm Movements: A
 Dynamical Systems Perspective. *Annual Review of Neuroscience*, *36*(1), 337-359.
 doi:10.1146/annurev-neuro-062111-150509
- Silberstein, M., & Chemero, A. (2013). Constraints on Localization and Decomposition as
 Explanatory Strategies in the Biological Sciences. *Philosophy of Science*, 80(5), 958-970.
 doi:10.1086/674533
- Stephens, R. G., Dunn, J. C., & Hayes, B. K. (2018). Are there two processes in reasoning? The
 dimensionality of inductive and deductive inferences. *Psychological Review*, *125*(2), 218 244. doi:10.1037/rev0000088
- Stepp, N., Chemero, A., & Turvey, M. T. (2011). Philosophy for the Rest of Cognitive Science.
 Topics in Cognitive Science, 3(2), 425-437. doi:<u>https://doi.org/10.1111/j.1756-</u>
 <u>8765.2011.01143.x</u>
- Sternberg, S. (1969). Memory scanning: mental processes revealed by reaction-time experiments.
 American Scientist, 57(4), 421-457. Retrieved from <u>http://www.jstor.org/stable/27828738</u>
- Strevens, M. (2013). No understanding without explanation. *Studies in history and philosophy of science part A*, 44(3), 510-515. Retrieved from
 http://www.sciencedirect.com/science/article/pii/S003936811200115X
- Sullivan, J. A. (2017). Coordinated pluralism as a means to facilitate integrative taxonomies of
 cognition. *Philosophical Explorations*, 20(2), 129-145. doi:10.1080/13869795.2017.1312497

- 1101 Tegmark, M. (2014). Our mathematical universe: My quest for the ultimate nature of reality:
 1102 Vintage.
- 1103 Thagard, P. (1978). The best explanation: criteria for theory choice. *Journal of philosophy*, 75, 76-92.
- 1104 Thagard, P. (1989). Explanatory coherence. *Behavioral and Brain Sciences*, *12*(3), 435-502.
- 1105 Thagard, P. (1992). Conceptual revolutions. Princeton: Princeton University Press.
- Thagard, P. (2012). *The cognitive science of science: Explanation, discovery, and conceptual change*.
 Cambridge: MIT Press.
- Thelen, E., Schöner, G., Scheier, C., & Smith, L. B. (2001). The dynamics of embodiment: A field
 theory of infant perseverative reaching. *Behavioral and Brain Sciences*, 24(1), 1-34.
- 1110 Turing, A. M. (1950). Computing Machinery and Intelligence. *Mind*, *LIX*(236), 433-460.
 1111 doi:10.1093/mind/LIX.236.433
- 1112 Ullman, S. (1979). The interpretation of visual motion. Cambridge: MIT Press.
- van Eck, D. (2018). Rethinking the explanatory power of dynamical models in cognitive science. *Philosophical Psychology*, *31*(8), 1131-1161. doi:10.1080/09515089.2018.1480755
- 1115 Van Hoeck, N., Watson, P. D., & Barbey, A. K. (2015). Cognitive neuroscience of human
 1116 counterfactual reasoning. *Frontiers in Human Neuroscience*, 9(420).
 1117 doi:10.3389/fnhum.2015.00420
- van Rooij, I., & Baggio, G. (2021). Theory Before the Test: How to Build High-Verisimilitude
 Explanatory Theories in Psychological Science. *Perspectives on Psychological Science*,
 16(4), 682-697. doi:10.1177/1745691620970604
- 1121 Venturelli, A. N. (2016). A Cautionary Contribution to the Philosophy of Explanation in the
 1122 Cognitive Neurosciences. *Minds and Machines*, 26(3), 259-285. doi:10.1007/s11023-016 1123 9395-0
- 1124 Verdejo, V. M. (2015). The systematicity challenge to anti-representational dynamicism. *Synthese*,
 1125 192(3), 701-722. doi:10.1007/s11229-014-0597-9
- 1126 Vernazzani, A. (2019). The structure of sensorimotor explanation. *Synthese*, *196*(11), 4527-4553.
 1127 doi:10.1007/s11229-017-1664-9
- Verreault-Julien, P. (2017). Non-causal understanding with economic models: the case of general
 equilibrium. *Journal of Economic Methodology*, 24(3), 297-317.
 doi:10.1080/1350178X.2017.1335424
- Wason, P. C., & Evans, J. S. B. T. (1974). Dual processes in reasoning? *Cognition*, 3(2), 141-154.
 doi:<u>https://doi.org/10.1016/0010-0277(74)90017-1</u>
- 1133 Watts, D. J., & Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 440-442. doi:10.1038/30918
- Weisberg, D. S., Keil, F. C., Goodstein, J., Rawson, E., & Gray, J. R. (2008). The Seductive Allure
 of Neuroscience Explanations. *Journal of Cognitive Neuroscience*, 20(3), 470-477.
 doi:10.1162/jocn.2008.20040
- Weiskopf, D. A. (2011). Models and mechanisms in psychological explanation. *Synthese*, 183(3),
 313-338. doi:10.1007/s11229-011-9958-9

- Wilkenfeld, D. A. (2013). Understanding as representation manipulability. *Synthese*, *190*(6), 9971016. doi:10.1007/s11229-011-0055-x
- Wilkenfeld, D. A. (2019). Understanding as compression. *Philosophical studies*, *176*(10), 2807-2831.
 doi:10.1007/s11098-018-1152-1
- Wilkenfeld, D. A. (2021). Objectually Understanding Informed Consent. *Analytic Philosophy*, 62(1),
 33-56.
- Williams, J. J., & Lombrozo, T. (2010). The Role of Explanation in Discovery and Generalization:
 Evidence From Category Learning. *Cognitive Science*, *34*(5), 776-806.
 doi:<u>https://doi.org/10.1111/j.1551-6709.2010.01113.x</u>
- Williams, J. J., Lombrozo, T., & Rehder, B. (2013). The hazards of explanation: Overgeneralization
 in the face of exceptions. *Journal of Experimental Psychology: General*, 142(4), 1006-1014.
 doi:10.1037/a0030996
- Woodward, J. (2003). *Making things happen: a theory of causal explanation*. New York: Oxford
 University Press.
- Woodward, J. (2013). Mechanistic Explanation: Its Scope and Limits. *Proceedings of the Aristotelian Society, Supplementary Volumes*, 87, 39-65. Retrieved from
 <u>http://www.jstor.org/stable/23482050</u>
- Zednik, C. (2011). The Nature of Dynamical Explanation. *Philosophy of Science*, 78(2), 238-263.
 doi:10.1086/659221
- Zipser, D., & Andersen, R. A. (1988). A back-propagation programmed network that simulates
 response properties of a subset of posterior parietal neurons. *Nature*, *331*(6158), 679-684.
 doi:10.1038/331679a0

1162