

What might interoceptive inference reveal about consciousness?

Abstract: The mainstream science of consciousness offers a few predominant views of how the brain gives rise to awareness. Chief among these are the Higher-Order-Thought Theory, Global Neuronal Workspace Theory, Integrated Information Theory, and hybrids thereof. In parallel, rapid development in predictive processing approaches have begun to outline concrete mechanisms by which interoceptive inference shapes selfhood, affect, and exteroceptive perception. Here, we consider these new approaches in terms of what they might offer our empirical, phenomenological, and philosophical understanding of consciousness and its neurobiological roots.

Keywords: Interoceptive Inference, Predictive Processing, Consciousness, Self-Inference, Global Neuronal Workspace, Higher Order Thought Theory, Active Inference

1	Table of Contents	
2	Table of Contents	2
3	Introduction	3
4	What is Consciousness?	3
5	What is Interoceptive Inference?	4
6	Predictive Higher Order Thought Theory (PHOTT)	8
7	GNWS and Interoceptive Inference	12
8	IIT and Interoceptive Inference	14
9	Active Inference, Interoception and Consciousness	18
10	Interoceptive Self-Inference: an integrated theory of consciousness?	20
11	Conclusion	22
12	References	24
13		
14		

15 Introduction

16 What is Consciousness?

17 If you have ever been under general anaesthesia, you surely remember the experience of waking
18 up. However, this awakening is different from the kind we do every morning, in that it is preceded
19 by a complete lack of subjective experience, a dark nothingness, without even the awareness of
20 time passing. This transition presents a clear insight into the two extremes of conscious experience.

21 While these strong contrasts delimit the borders of consciousness, you might also consider
22 the phenomenological properties which reveal themselves upon further reflection. Foremost is the
23 unique “mineness” of any conscious experience. In the transition from sleep to wakefulness, there
24 seem to be distinct properties of ownership and agency. Whereas the infinite void of sleep belongs
25 to no one, even before opening my eyes there is a distinct sense in which experience is happening
26 to someone. In phenomenological terms, we can think about this as the minimal, pre-reflexive
27 conditions about which my experiences are uniquely my own (Gallagher, 2000). Consciousness
28 then is something which happens to a sentient subject, which is lived through as the embodied
29 point of view of those seemingly ineffable subjective properties.

30 A sufficient theory of consciousness then, will deal with each of these properties in turn.
31 What distinguishes conscious states from non-conscious ones? How does selfhood and agency
32 influence these properties? Which sorts of mechanisms give rise to both the phenomenological
33 contents of consciousness, and determine which sorts of states become accessible to conscious
34 thought? How might the body, or emotions, interact with these properties of consciousness?

35 Answering these questions is no easy task. Certainly, most who study consciousness have
36 heard the joke that there are as many theories of consciousness as there are consciousness theorists.
37 Our goal here is not to provide a comprehensive predictive processing or active inference theory
38 of consciousness, of which there are already a rapidly growing number (for reviews, see Hohwy
39 & Seth, 2020; A. K. Seth & Hohwy, 2020; Whyte, 2019; Whyte & Smith, 2021). Rather, we aim
40 to illustrate how the notion of interoceptive inference and related concepts might inform the
41 theoretical and empirical science of consciousness, by generating alternative process theories that
42 can then be subject to empirical evaluation.

43 Current mainstream approaches to consciousness can be largely divided into several
44 camps, though the boundaries are fuzzy and hybrid theories abound. Writ large, these include the
45 Global Neuronal Workspace Theory (GNWS), Higher-Order Thought Theory (HOTT), and the
46 Integrated Information Theory (ITT). These theories share some key properties, but also differ
47 substantially in terms of the types of phenomena they seek to explain and the mechanisms they
48 appeal to in doing so. In what follows, we will discuss some of the more obvious places in which
49 predictive processing and interoceptive inference theories tie in with these approaches. Here, we
50 summarize key concepts from some of the leading theories of consciousness and discuss how
51 interoceptive inference might fit into them and inform future theoretical and empirical directions.
52 Our main goals here are the following; first, to accurately and concisely review several of the most

53 popular theories of consciousness, namely HOTT, GWS, IIT and active inference accounts. We
54 then aim to describe the emerging concept of interoceptive inference, and finally we explore the
55 potential of interoceptive inference to integrate with each of the theories, and how it might
56 illuminate future research directions.

57 What is Interoceptive Inference?

58 First, however, we must introduce the standard set pieces of predictive processing and
59 interoceptive inference. Predictive processing can be described as a set of theories which aim to
60 understand how expectations – both neural and psychological – shape, constrain, and ultimately
61 define the mind. These theories have deep roots in cybernetics, information processing, and
62 seminal prospective control models emerging from early 1960s motor and activity theory. A key
63 feature of predictive processing is the basic notion that biological information processing occurs
64 primarily via the minimization of (information theoretic) surprise, such that the nervous system
65 can be understood as a hierarchy of top-down predictions and bottom-up prediction errors. Whilst
66 most early theories extrapolated this basic scheme to explain restricted phenomena such as
67 prospective motor control and the sense of agency (Sperry, 1950; Synofzik et al., 2008; von
68 Helmholtz, 1925), in recent years these approaches have exploded with a myriad of conceptual,
69 computational, and empirical work.

70 An in-depth review of the scope of predictive processing is beyond this current article. For
71 the unfamiliar reader, we here recall the basic principles, but for a more thorough treatment
72 numerous recent reviews exist, both of the general computational and theoretical principles (Bastos
73 et al., 2012; Clark, 2013; K. Friston, 2009, 2018b; Hohwy, 2013), and their relationship with
74 notions of embodiment and selfhood (Petzschner et al., 2021; A. Seth & Critchley, 2013; A. K.
75 Seth & Friston, 2016).

76 In summary, these approaches surmise that the brain, much like a Russian Matryoshka or
77 nesting doll, comprises an interlocking hierarchical web, with each unit or level of this web
78 predicting the output of the lower level. At the outermost layer of this hierarchical ‘brain web’ one
79 finds the sensory epithelium and motor apparatus of the agent – that is, the means by which the
80 agent takes in information about the world external to itself, and acts upon those sensory inputs to
81 alter the world. As one moves from these outermost layers, venturing deeper into the nervous
82 system, neuronal populations encode or invert a model of its inputs¹. This generative model
83 comprises three key components: a prediction (e.g., of a hierarchically lower expectation), a
84 prediction error (e.g., encoding the difference between the expectation and its prediction), and the
85 precision of each of these signals (e.g., encoding their predictability). This simple motif is
86 replicated from the lowest, most basic neural representations of first order neurons predicting the

¹ invert here is using the technical (Bayesian sense) it refers to the inverse mapping between consequences and causes afforded by a generative model where causes generate consequences. In short, inverting a generative model means inferring the (hidden) causes of (observable) consequences.

87 activity of sensory effectors, to the highest order, most polymodal representations encoding
88 concepts, selfhood, and preferences.

89 Early predictive processing theories largely appealed to this motif of prediction error
90 minimization (PEM) to explain phenomena such as visual perception (Rao & Ballard, 1999), motor
91 control (K. Friston, 2011), agency (Synofzik et al., 2008), or social cognitive meta-representation
92 (Kilner et al., 2007; Koster-Hale & Saxe, 2013; Tamir & Thornton, 2018). In contrast, the new
93 “radical predictive processing” wave embraces the unifying nature of the predictive brain in an
94 attempt to explain how all aspects of information processing and behaviour emerge from the
95 integrated hierarchical flow of predictions, prediction errors and their precision (Allen & Friston,
96 2018; Clark, 2015). Within this framework then, we can consider both the specific hierarchical
97 processing of interoceptive sensations (Allen, 2020; A. K. Seth, 2013a), and the broader
98 implications of embodied, affective inference with respect to our understanding of consciousness.

99 Interoception is generally used to refer to the sensation, perception, and metacognition of
100 the visceral cycles which govern an agent’s homeostasis, allostasis, and ultimately its survival
101 (Barrett & Simmons, 2015; Sherrington, 1952; Vaitl, 1996). This includes, on the ascending side,
102 the sensory information conveying heartbeats, respiration, and the activity of the stomach and gut
103 to the brain – literally, gut feelings. On the descending side, interoception denotes the visceromotor
104 signals and allostatic reflex arcs by which agents maintain their homeostasis in the face of
105 environmental challenges. Interoceptive processes are thus those which enable an agent to monitor
106 and control the bodily states that are necessary to maintain the balance between energy expenditure
107 and consumption.

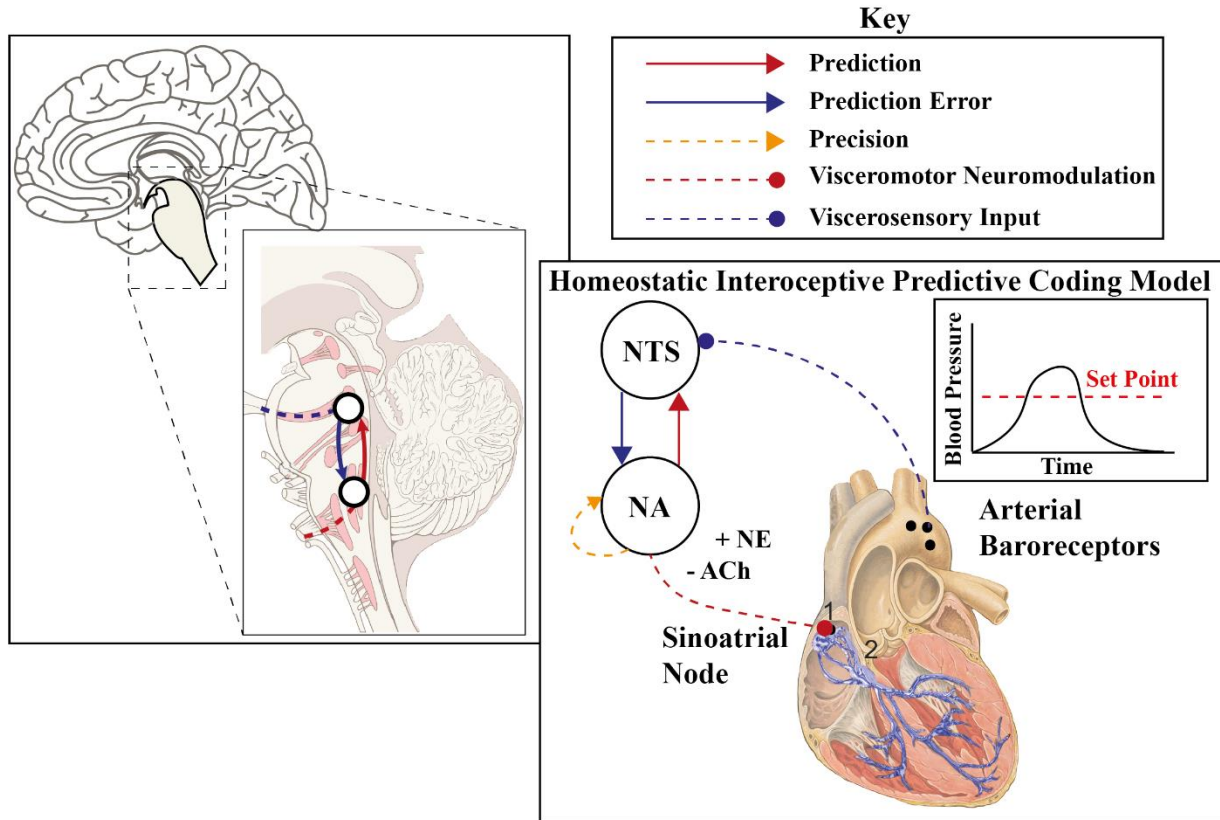
108 We can further demarcate interoceptive processes into those which directly subserve
109 *homeostasis*, that is the maintenance of a steady state defined by specific metabolic set-points, and
110 *allostasis*, the proactive control of the body – and environment – to resolve homeostatic needs
111 before they arise (Barrett et al., 2016; Kleckner et al., 2017; Sterling & Eyer, 1988). For example,
112 biological necessity dictates that body temperature, blood oxygenation, and blood glucose level
113 are all maintained within a restrictive range of values. Any sustained deviation from these values
114 is likely to negatively impact an organism’s survival, whether through the direct inducement of
115 cellular death, or by the slow attrition of metabolic surplus through starvation. If oxygen is too
116 low, or temperature too high, the brain can directly engage adaptive physiological reflexes,
117 maintaining homeostasis by increasing respiratory frequency or decreasing systolic blood
118 pressure.

119 These simple sensory-motor reflex arcs, illustrated in Figure 1, can be readily understood
120 by appeal to predictive mechanisms not unlike that of a common household thermostat. That is to
121 say, a low-level spinal, thalamic, or brainstem circuit is generally sufficient to encode the set-point
122 as a prior expectation on the heart-rate, respiratory frequency, or blood pressure. As in afferent
123 control theory, this problem reduces to one of increasing or decreasing the descending
124 visceromotor predictions to minimize any sensory prediction error that occurs: c.f., the equilibrium
125 point hypothesis in motor control (A. G. Feldman, 2009) and related perceptual control theories
126 (Mansell, 2011). One can thus easily envision simple predictive engrams, which monitor visceral

Interoceptive Inference and Consciousness

127 inputs and adjust bodily states as needed to maintain the overall integrity of the system. By
128 comparing the re-afferent sensory inputs to the expected change induced by each top-down
129 prediction, the system can meet whatever thermoregulatory, metabolic, or other homeostatic
130 demands are needed, with relatively little need for higher order cognition.

131
132
133
134



135
136 **Figure 1: Simplified Homeostatic Control via Interoceptive Inference.** This simplified schematic illustrates an
137 example of low-level interoceptive predictive coding in the cardiac domain. Here, a simplified two-node control loop
138 maintains a homeostatic set-point by minimizing the error between afferent cardiac sensory inputs and descending
139 neuromodulatory efferent control. In this example, blood pressure and heart-rate are controlled by a cardiac
140 comparator circuit circumscribed in the primary medulla of the brainstem. Arterial baroreceptors located in the aorta
141 and carotid artery increase their firing rate whenever blood pressure rises above a homeostatic set-point. This firing
142 is relayed via the cranial nerves to the nucleus tractus solitarius (NTS). The NTS acts as a comparator, computing the
143 difference between descending blood pressure predictions and these incoming signals. The difference, or prediction
144 error, is relayed upwards to the nucleus ambiguus (NA), which regulates heart rate via descending cholinergic
145 neuromodulation, triggering the sinoatrial node (SA) to reduce cardiac frequency. An efferent copy (i.e., a descending
146 cardiac prediction) is sent downwards to the NTS, and the comparative loop continues until blood pressure falls below
147 the homeostatic set-point. The relative strength of top-down and bottom-up signals (i.e., their precision) is regulated
148 via neuromodulatory gain control, depicted as self-connections in orange. For illustration purposes, the underlying
149 cardiac neurophysiology has been simplified, leaving out for example the perfused excitatory effects of noradrenaline.
150

151 In contrast, allostatic processes are needed whenever the environment or body can no longer
152 maintain these set-points through simple, internal reflex actions alone. For example, if I
153 consistently fail to meet my energy needs, the body will begin to consume itself. Here, merely
154 maintaining homeostasis is insufficient for survival – the agent must identify the external, hidden
155 causes which are causing the increased allostatic load. For example, the environment may no
156 longer contain sufficient resources, in which case the agent should deploy exploratory cognitive
157 mechanisms to find greener pastures. Similarly, if an environment becomes overly threatening
158 (i.e., if the long-term volatility of threats increases), merely increasing or decreasing the heart-rate
159 is no longer sufficient. Instead, I should engage more complex fight or flight routines, to remove
160 the immediate threats and make the situation more amenable to my survival.

161 Interoceptive inference in the context of allostasis can thus be viewed as operating at a
162 level once (or thrice) removed from that of basic homeostasis. Whereas interoceptive inference at
163 the first order might merely involve the regulation of viscerosensory and visceromotor prediction
164 errors, allostatic interoceptive inference requires the agent to link these low-level variables to
165 contextual ones operating at fundamentally longer timescales. As such, interoceptive inference at
166 this level naturally links to the representation of selfhood, valence, and other metacognitive
167 concepts linking the agent’s current homeostatic state to the overall volatility of its environment
168 and conspecifics (Barrett et al., 2016; Petzschner et al., 2017, 2021).

169 At a broader level still, we can consider the phylogenetic and ontogenetic role that
170 interoceptive processes play in the overall structure and organization of the nervous system. One
171 standout example of this is found in the Free Energy Principle (FEP), a normative biological theory
172 which posits specific foundational and information theoretic constraints on specific biological
173 process theories (K. Friston, 2009). The FEP emphasizes that at the very basis of any biological
174 agent is the self-organized maintenance of its own existence (K. Friston, 2013). In this sense, the
175 very structure of the nervous system can be seen as entailing a generative model², which ensures
176 the agent will engage in both homeostatic and allostatic processes. Under the FEP then, the body
177 (both visceral and somato-morphic) are understood as a kind of “first prior” (Allen & Tsakiris,
178 2018; Ciaunica et al., 2021), which shapes the evolutionary refinement of the predictive mind.
179 Through this lens, the interoceptive hierarchy plays a special role not only in maintaining an agents'
180 survival, but in determining the salience of every action and ensuing belief updating, and
181 ultimately value itself is understood as whatever maximizes the evidence for the agents' model of
182 a survivable world (c.f., “the self-evidencing brain”) (Hohwy, 2016).

183 What then can these set pieces about the brain tell us about consciousness? To start, any
184 predictive processing theory will obviously posit a central role for expectations and predictions in
185 the genesis and contents of consciousness. If the mind is primarily concerned with the
186 representation of future events (i.e., the consequences of action), then it seems likely that

² ‘Entail’ is used carefully here to acknowledge that the generative model is a mathematical construct, not something that is physically realized: neuronal processes can be understood as minimizing free energy that is a function of a generative model; however, neuronal dynamics that are realized reflect *free energy gradients* (that can be cast as a prediction error), not the free energy *per se*.

187 consciousness is also predominantly prospective. But should a theory of consciousness posit that
188 specific, higher-order modules generate our subjective experience, or rather that it emerges from
189 the collective prediction error minimizing activity of the organism? Similarly, it can be assumed
190 that most predictive processing theories of consciousness should posit a central role in the
191 encoding and modulation of precision – in determining which particular predictions become
192 conscious, and in terms of how conscious predictions should influence affective, metacognitive,
193 and self-related phenomena. That is to say, a basic predictive processing theory of consciousness
194 is likely to ascribe some facets of both *access* and *phenomenal* consciousness (Block, 1995) to
195 error minimizing predictions, and the precision of signals which ensure one particular hypothesis
196 dictates the contents of consciousness versus another.

197 Must we appeal at all to interoceptive processes in a theory of consciousness? As we shall
198 see, this depends largely on the overarching theory of consciousness developed, i.e., which
199 conscious phenomena are the target of explanation. Certainly a general PEM-based theory of
200 consciousness would ascribe our bodily self-consciousness to the hierarchical minimization of
201 homeostatic and allostatic prediction errors (Ainley et al., 2016). E.g., my consciousness of my
202 heart rate or respiration could be argued to be a product of the viscerosensory and visceromotor
203 prediction errors and precision signals which drive Bayesian belief updating. In this sense, such
204 interoceptive sensations are likely to dominate my awareness, whenever these systems give off
205 prediction errors, whose precision may need updating³. Yet such a theory would not posit anything
206 particularly unique about interoceptive inference, casting it as just another parcel of the
207 hierarchical organism which gives rise to various bodily aspects of consciousness. Alternatively,
208 one could develop an FEP or similarly radical predictive-processing theory of consciousness,
209 wherein interoceptive inference may fundamentally underpin access and/or phenomenal
210 consciousness. To consider these different possibilities, we now review predominant theories of
211 consciousness in light of interoceptive inference.

212 Predictive Higher Order Thought Theory (PHOTT)

213
214 Higher order-thought theories (HOTT) stem originally from the analytic philosophy of mind
215 (Carruthers, 2007; Carruthers & Gennaro, 2020; Rosenthal, 2006), yet have also found substantive
216 purchase in the empirical science of consciousness (Brown et al., 2019; H. Lau & Rosenthal, 2011;
217 LeDoux & Brown, 2017). In essence, HOT theories argue that properties of conscious experience
218 arise from the relationship between mental states and higher-order representations of these states
219 (Rosenthal 2005). Critically, this implies that a first-order representation by itself is not part of
220 conscious content, unless it is accompanied by another (higher-order) process that is reflecting on

³ The updating of the precision of prediction errors is generally read as sensory attention or attenuation Brown, H., R. A. Adams, I. Parees, M. Edwards and K. Friston (2013). "Active inference, sensory attenuation and illusions." *Cogn Process* **14**(4): 411-427. This speaks to an intimate link between conscious (interception-pointing) inference and attentional selection – or sensory attenuation.

221 its content. In this sense, HOTT stipulates that an agent can be conscious of some representation
222 *X if and only if* the agent possesses a higher order meta-representation of *X*. This approach is based
223 on strong assumptions about the links between phenomenal and access consciousness: according
224 to HOTT, conscious states are by definition those that the agent is aware of.

225 Empirically speaking, HOTT is often associated with metacognitive approaches to
226 modelling consciousness, such as the popular signal-detection theoretic (SDT) framework (Brown
227 et al., 2019; Fleming & Lau, 2014; Ko & Lau, 2012; H. C. Lau, 2007; H. Lau & Rosenthal, 2011).
228 Here, for a conscious state to be labeled as such, an experimental subject must not only exhibit
229 above chance accuracy for detecting some stimulus, but also show explicit conscious awareness
230 of their own accuracy, typically measured via subjective confidence or awareness ratings. Now the
231 metric of consciousness is not just whether a subject can reliably discriminate or detect some input,
232 but whether the subject possesses an accurate meta-representation of their own sensory process
233 (i.e., there should also be a high correlation of confidence and accuracy). Neurobiologically,
234 HOTT proponents frequently argue that the prefrontal cortex plays a crucial role in this
235 metacognitive re-representation of first-order perceptual contents, and as such is sometimes said
236 to be a necessary and sufficient neural correlate of consciousness (NCC).

237 What then might a “predictive higher-order thought theory” (PHOTT) look like? To our
238 knowledge no theorists have yet directly developed a PHOTT, and a full derivation is beyond the
239 scope (and expertise) of the present article. However, we here briefly sketch some constitutive
240 components of a potential PHOTT, in the hopes of illuminating how interoceptive inference might
241 contribute to such a theory, and in guiding future theoretical, empirical, and computational work.

242 While thus far no explicit PHOTT theory of consciousness has been proposed, the close
243 alignment of these approaches to empirical and computational metacognition research provides
244 some clear starting points. Metacognition, i.e., the meta-representation and control of first order
245 cognitive or perceptual processes, is typically viewed through a decision-theoretic framework in
246 which the agent must monitor the signal and noise distributions underlying first-order perceptual
247 performance, in order to arrive at a representation of the overall probability that one is making
248 correct responses. Fleming (2020) suggests starting with metacognitive reports of awareness; after
249 all, we can only be aware of another’s conscious state through their reports, through language. In
250 the Higher Order State Space (HOSS) model, awareness corresponds to inference on the generative
251 model of the perceptual content, and can be represented as an additional hierarchical state that
252 signals whether perceptual content is present or absent in lower levels. In this case, the higher order
253 thought is cast as a posterior belief over the lower-order contents of consciousness.

254 Several theorists have proposed Bayesian or predictive extensions of these basic models
255 (Fleming & Daw, 2016; H. C. Lau, 2007; Yeung & Summerfield, 2012), where typically the
256 second-order model is seen as integrating either the contents or the precision of lower-order
257 representations (e.g., the confidence associated with a prediction error encoding a visual input)
258 with high-level “self-priors” describing one’s efficacy or overall ability within that cognitive
259 domain. Thus, a basic Bayesian view of metacognition (and meta-representation more generally)

260 posits an extended cognitive hierarchy in which low-level precision signals are read-out and
261 integrated according to some higher-order self-model.

262 This raises some immediate set-pieces and questions for a PHOTT model of consciousness.
263 In the philosophical literature the exact nature of the meta-representation needed to render a first
264 order representation conscious has been the subject of intense debate. For example, opposing
265 philosophical camps argue that a HOT must be conceptual in nature to render phenomenal
266 consciousness, versus “higher order perception” (HOP) theorists who posit a kind of “inner sense”
267 theory, which maintains that HOTs need not be conceptual in nature.

268 Returning to the predictive brain, we find multiple possible candidates for HOTs or HOPs,
269 depending on what particular process theory one works within. For example, in more modular or
270 comparator-based approaches to predictive processing, one could posit the existence of an explicit
271 metacognition module which monitors first-order perceptual representations in order to form an
272 explicit, conceptual HOT encoding the probability that these are correct (as opposed to illusory)
273 percepts. In this sense, predictive higher-order-thoughts (PHOTS) would be ascribed to the higher-
274 order, content-based predictions originating from deep within the brain’s hierarchy, encoding
275 relational properties between conscious contents (e.g., the connection between the sensory features
276 encoding a lover’s face and the warm affective association therein), or as in the Bayesian
277 metacognitive modules described before, simply encoding the prior probability that a percept is
278 correct given some conceptual self-knowledge and the ongoing pattern of lower-order perceptual
279 prediction errors.

280 Alternatively, one could argue for a PHOTT (or perhaps a PHOPT) in which the contents
281 of first-order prediction are largely irrelevant to whether a percept becomes conscious or not, and
282 instead emphasize that PHOTs are fundamentally concerned with meta-representing the *precision*
283 of lower-order contents. This aligns both with extant Bayesian theories of metacognition, which
284 emphasize that subjective awareness arises from a posterior estimate of precision, and with the
285 intuitive notion that precision is itself fundamentally a second-order statistic (that is, a meta-
286 representation) of first-order predictive processes. In this case, a precision-focused PHOT would
287 likely emphasize the role of higher-order neural modules in extracting and re-representing the
288 precision (but not the contents) of lower-order predictions, and conscious states would be those
289 associated with the greatest a posteriori precision.

290 Clearly, these examples are meant to serve as high level outlines illustrating how the set-
291 pieces and explanatory concepts present in predictive processing can be circumscribed within a
292 HOTT of consciousness. Much work remains to be done extrapolating from these basic ideas to a
293 rigorous overall theory. We anticipate that along the way, difficult questions will need to be
294 addressed, concerning for example whether PHOTS are fundamentally concerned with contentful
295 meta-representation, or only with representing the confidence or predictability of first-order
296 processes. One interesting question which emerges immediately, for example, is whether any
297 precision signal could be seen as a sufficient higher-order meta-representation, or whether only
298 higher-order *expected precision* signals would qualify. What we mean is that, according to radical

299 predictive processing theories (Clark, 2015), precision signals can be found at all levels of the
300 central nervous system (Allen & Tsakiris, 2018; Bruineberg & Rietveld, 2014).

301 At each level of the brain's canonical microcircuitry then, there is a kind of meta-
302 representation encoding the precision of prediction errors arising at that level, and these local
303 precision signals govern the overall flow of contents through the cortical hierarchy. Are these low-
304 level meta-representations sufficient for a content to become conscious? If so, it would appear then
305 that a PHOT theory of consciousness may help to unify recurrent neural processing and HOTT
306 approaches (Lamme, 2006; Lamme & Roelfsema, 2000), as phenomenal consciousness would
307 emerge from the interaction of local recurrent connections and their associated precision weighting
308 low-level perceptual circuits. In contrast, if it is the explicit representation of expected precision
309 (i.e., top-down, typically poly-modal predictions of future changes in lower order precision) that
310 renders a lower state conscious or not, then the resulting PHOT would likely ascribe
311 neuromodulatory circuits and prefrontal modules as fundamental for determining consciousness
312 (Boly et al., 2017; Odegaard et al., 2017).

313 How does interoception fit into the PHOTT framework? One option is that interoceptive
314 information, just like visual input, is another source of lower order perceptual input, which can be
315 integrated with other information and reflected upon by higher order processes to become a subject
316 of conscious experience. In this sense then, PHOTs predicting either higher-order interoceptive
317 contents (e.g., the association between multiple viscerosensory systems and affect or value) would
318 largely determine whether one is conscious or not of any given interoceptive sensation. In this
319 sense, interoception would not play any special role in a PHOT theory of consciousness, other than
320 offering another channel of perceptual contents which may be configured within any other higher
321 order thoughts or percepts.

322 Alternatively, if the preferred PHOTT emphasizes the role of meta-representations
323 encoding expected precision, then interoceptive processes may play a more constitutive role in
324 determining either phenomenal or access consciousness. Generally speaking, the optimization of
325 expected precision has been proffered as a unifying mechanism by which salience, attention, and
326 high-level self-control emerge (H. Feldman & Friston, 2010; Kanai et al., 2015; Moran et al., 2013;
327 Parr & Friston, 2017, 2019). Furthermore, the very capacity to supply low levels of hierarchical
328 inference with predictions of precision or predictability has been proposed as a necessary condition
329 for qualitative experience; in the sense of precluding phenomenal transparency (Limanowski,
330 2017; Limanowski & Blankenburg, 2013; Limanowski & Friston, 2018).

331 This approach views bottom-up and top-down attention as emergent properties of minimizing
332 "precision-prediction errors", such that the top-down control of expected precision can selectively
333 enhance or inhibit lower-order percepts. Interoceptive prediction errors and precision thereof are
334 here thought to play a unique role in determining what is salient for an agent in any given context,
335 such that unexpected challenges to homeostasis or allostasis essentially govern the innate value of
336 different outcomes. Computational and conceptual models have expanded on this view to describe
337 a process of metacognitive and interoceptive self-inference, in which the a priori expected
338 precision afforded the homeostatic and allostatic fluctuations is always higher than that of over

339 sensory-motor channels (Allen, 2020; Allen et al., 2019, 2020). As fluctuations in, for example,
340 blood temperature or arterial pulsation, can directly modulate the noise (i.e., inverse precision) of
341 neuronal circuits in a global fashion (Chow et al., 2020), then the representation of expected
342 precision is argued to both sample directly from the precision of interoceptive prediction errors,
343 and to utilize descending visceromotor control as a means of optimizing sensory precision.

344 In PHOTT terms then, this could be taken as an argument that visceral prediction and
345 precision signals play an especially important role in the meta-representation of first-order
346 perceptual contents, such that either their subjective salience is largely governed by higher-order
347 thoughts about the interaction between the visceral body and the exteroceptive sensorium. In this
348 sense then, both the “shape” or “contours” of phenomenal consciousness, and the likelihood that
349 a percept becomes conscious (i.e., access consciousness) may depend in part on the top-down
350 meta-representation of expected (interoceptive) precision. Such a process theory would then show
351 some alignment between the PHOTT approach, and recent theoretical proposals suggesting that
352 interoceptive signals may play a fundamental role in shaping the “mineness” or subjective quality
353 of conscious experiences (Ainley et al., 2012, 2016; Azzalini et al., 2019; Fotopoulou & Tsakiris,
354 2017; A. K. Seth & Tsakiris, 2018).

355 GNWS and Interoceptive Inference

356 The Global Workspace theories (Baars, 1988; Baars & Franklin, 2007; Dehaene & Changeux,
357 2011) originate from the idea that consciousness arises from processing of information in the brain,
358 and the way in which specific information is selected and broadcast across the brain in order to
359 generate a coherent representation. Here, the brain is composed of a set of specialized, local
360 cortical processing units, which are richly interconnected by excitatory pyramidal neurons
361 spanning between frontal and parietal regions. A piece of information, represented in one or several
362 of the processing units, can cross a threshold and be selected for broadcasting (i.e., amplification)
363 in the process of ‘ignition’, whereby it is simultaneously made available to all processing units.
364 For example, when a bird perches nearby and chirps, my attention is drawn to the sound, and I will
365 gaze around to find the source. The perceptual inputs associated with the bird are carried up and
366 processed, and as they enter the global workspace and become ‘globally available’ as a part of
367 consciousness, such that, they, along with the idea of the bird and the feeling the moment is
368 associated with in my body, are broadcast to various brain systems.

369 These may include memory allowing me to remember the moment, motor action, or higher
370 cognitive systems which enable me to make decisions and talk about my experience. Crucially,
371 most information that is available to and processed by the brain need not enter the global
372 workspace, here consciousness is about how and which information is selected for global
373 processing and awareness. The Global Neuronal Workspace (GNW) theory (Dehaene et al., 2006;
374 Dehaene & Changeux, 2011) specifies that information which does become available to the global
375 workspace (GW) then recruits brain networks extending over frontal and parietal regions which
376 can integrate the dispersed sources of information into a coherent conscious phenomenon. Thus,

377 the prefrontal cortex plays a central role within GNW as in HOT theories, yet they differ in what
378 functions they ascribe to it (Mashour et al., 2020); in HOT the higher-order metacognitive
379 processes representing first-order states are what constitute consciousness, so if they are in the
380 PFC, this region becomes a source of consciousness. In GNW meanwhile, conscious states emerge
381 by the broadcasting of information across systems, which can happen due to long range
382 connections between PFC, other fronto-parietal regions comprising the GW. We emphasize that
383 HOTT and GNW are not mutually exclusive, and in fact, several works aim to bridge and unify
384 these theories (Dehaene et al., 2017; Graziano et al., 2019). The Attentional Schema Theory for
385 example, merges GNW and HOT by proposing that attention amplifies signals so that they may
386 reach ignition, and that there is a higher order representation of the GW which represents the
387 dynamics and implications of having a GW, which is what gives rise to phenomenological
388 consciousness.

389 Unifying approaches are also building active inference-based models within the GNW
390 framework. The predictive global neuronal workspace (PGNW) (Hohwy, 2013; Whyte, 2019;
391 Whyte & Smith, 2021) combines Bayesian active inference with experimentally corroborated
392 components of the GNW. The PGNW enables us to examine one of the core questions arising from
393 the GNW theories: what determines the ignition threshold? Within predictive processing, the
394 information that crosses the threshold to reach ignition is that which best accommodates PEs
395 throughout the hierarchy, so that the best-fitting (PEM) model of the world is selected and
396 broadcast across systems (Friston, Breakspear et al. 2012). Ignition then represents the point at
397 which an evidence accumulation process has reached the threshold where it becomes the most
398 likely explanation of the world (i.e., the current ascending input). The PGNW therefore represents
399 ignition as an inferential process. As in the active inference framework below, ignition here
400 requires sufficient temporal thickness to coordinate and contextualize lower levels of processing
401 (Whyte & Smith, 2021). In order to be able to speak of my experience of the chirping bird, I need
402 a representation that is maintained for some period of time and that extends back in time to include
403 me observing the bird.

404 According to the standard GNW account the anterior insula, a key hub processing visceral
405 information and involved in interoceptive awareness (Critchley et al., 2004; Evrard, 2019), selects
406 and prioritizes information prior to possible amplification by the GW (Michel, 2017). Another
407 theory in the same spirit (Chanes & Barrett, 2016), presents the limbic cortex (including the
408 anterior insula, anterior cingulate cortex, among other areas) as the ‘limbic workspace’ in light of
409 the rich bi-directional connections between these areas and lower levels of processing. In this view,
410 cortical lamination is a distinguishing feature, so that predictions move up from less to more
411 laminated areas while PEs move down in the opposite direction.

412 Within the PGNW view, interoception is a perceptual system (or set of systems), sensing the
413 internal states and rhythms occurring in the body, and information from it can independently or
414 together with congruent information from other systems, be broadcast by the GNW. For example,
415 I may become aware of a sudden stomach cramp, which incites me to think about what I have
416 eaten earlier in the day. However, recent evidence proposes that interoception might also play a

417 modulating role on other systems (Chanes & Barrett, 2016; Michel, 2017), whereby interoceptive
418 prediction errors and/or precisions affect the likelihood that other modules are brought into the
419 GW, driving ignition itself through the modulation of salience.

420 It has been suggested that the brain maintains a self-model representing the status of the
421 body, which is continuously updated to fit ascending interoceptive input by changing interoceptive
422 PEs (Barrett & Simmons, 2015). Generally, in these accounts, what achieves ignition can be
423 understood as the relationship between the expected (top-down) and sensory (bottom-up)
424 precision, where when the self-model increases the precision of lower order modules, they become
425 better fitting models of the world and are more likely to reach ignition. Further work has proposed
426 that interoception may play a crucial role within the self-model, by either conditioning expected
427 precision (Allen et al., 2019, 2020), or by modulating the degree to which lower-order
428 representations are interpreted as related to the sense of self (Apps & Tsakiris, 2014; Azzalini et
429 al., 2019; Babo-Rebelo et al., 2016; Limanowski & Blankenburg, 2013; Quigley et al., 2021; A.
430 K. Seth, 2013b).

431 Thus, as we saw in the previous section, depending on the exact predictive process theory
432 one motivates, interoception may act simply as one of many modules within the GNWS, or it may
433 play a more foundational role, either by guiding the top-down selection of modules into the WS
434 by the self-model, or by enhancing the gain or precision associated with lower-order, non-
435 interoceptive modules as to alter their probability of promotion into the WS. We therefor propose
436 that future experimental and computational work will likely benefit from modelling how
437 interoceptive processes interact with conscious processing of stimuli, and the proposed
438 neurophysiological signatures of ignition, such as the P300 component, to ultimately understand
439 whether interoceptive prediction errors or their precision alter the process of ignition and the
440 overall topology of the GNWS.

441 IIT and Interoceptive Inference

442 The Integrated Information Theory (IIT) (Oizumi et al., 2014; Tononi et al., 2016) of
443 consciousness is an attempt at a formal method for mathematically describing the conscious
444 experience of any given system, agnostic as to all but the causal structure of its substrate. The
445 theory focuses on making an intrinsic description of the system, that is, how the system is to itself,
446 opposed to an extrinsic description from the perspective of an outside observer. The IIT takes as
447 starting point five axioms for what constitutes any phenomenological experience. From that, five
448 criteria are derived which must be met in order for a physical system to support conscious
449 experience: 1) the system must exist, that is, exert and be subject to causal power; and it must do
450 so over itself in a way that is 2) structured of component elements; 3) informative i.e.,
451 distinguishable from other causal states; 4) integrated or unitary as a whole, and irreducible to
452 independent subsets; and 5) exclusive or definite, specifying its own borders.

453 To measure the degree to which a system fulfils these criteria, a measure of *integrated*
454 *conceptual information* is used, denoted as Φ , which measures the degree to which the system

455 exerts causal power over itself in a way that is irreducible to the activity of its components. The
456 conscious parts of a system are called complexes and are those parts of the system that specify the
457 highest Φ without overlapping with one other. The axiom that complexes cannot overlap also
458 means that smaller complexes are not conscious, even if Φ is larger than zero, as long as they are
459 contained within a larger complex with a higher Φ . Conversely, a large complex is also not
460 conscious if there are smaller complexes within it with a higher Φ . This leads to predictions of
461 consciousness in the brain being situated in areas with more integrated connections, currently
462 thought to be an temporo-parietal-occipital hot-zone in the posterior cortex (Koch et al., 2016).
463 This excludes more feed-forward networks like the cerebellum, explaining why this structure does
464 not obviously contribute to consciousness despite its large number of neurons (Tononi & Koch,
465 2015). The exclusivity axiom also means that experience only happens at one spatial and temporal
466 scale of organization, namely at the level at which Φ is highest (Hoel et al., 2016).

467 The conscious experience of a complex involves concepts, which are causal mechanisms
468 within the complex that specify irreducible cause-effect repertoires. All the concepts together form
469 a concept structure, which can be interpreted as a geometric shape in a multidimensional concept
470 space. The concept structure of a complex is thought to reflect the *content* of consciousness, while
471 the size of Φ reflects the *amount* of consciousness as a whole. Importantly, the concept structure
472 not only depends on the current state of the complex, but also on the other possible states it could
473 take, since it is defined by how the system causally constraints itself. This allows for the possibility
474 for negative concepts, that is, the absence of some state (e.g., not-red), and that conscious
475 experience is also enriched by the increase of more possible states (e.g., seeing green includes not-
476 red, not-blue, etc.).

477 The IIT and predictive processing theories of brain function and of consciousness take
478 quite different starting points in a range of respects: IIT is concerned with understanding how
479 systems in general relate to themselves, while predictive processing addresses how the brain,
480 specifically, relates functionally to the surrounding environment, including the body. The former
481 begins entirely in describing phenomenology to identify compatible types of physical systems,
482 while the latter largely takes the opposite direction and starts with what is required for the physical
483 brain in order to self-organize and maintain itself, going from there to describe phenomenology.
484 This makes it challenging to combine the two approaches, a project that is far beyond the scope of
485 this paper – but see (K. J. Friston et al., 2020) for a discussion, in terms of the information geometry
486 of active inference.

487 It might be worth briefly speculating, however, what predictive processing accounts might
488 be able to offer IIT to inform the broader discussion of interoceptive processing and consciousness.
489 One notion is that hierarchical, precision-weighted prediction error belief updating schemes might
490 provide (neuronal) structures that result in high levels of integration, a suggestion that might
491 potentially be investigated by calculating Φ of canonical neuronal schemas from predictive
492 processing and comparing it to other proposed schemas. The prediction error minimization loops
493 in PEM theories are certainly more complex and integrated than the zero Φ feedforward networks
494 in, say, artificial neural networks. Zooming out, one might also ask if the overall structure of the

495 brain relates to the level of integration; does, for example, the presence of a self-model in the brain
496 somehow constitute or allow for higher levels of integration? One could certainly imagine that the
497 part of the brain that constitutes the ‘self-as-hypothesis’ might be highly integrated, given that it
498 has to coordinate impressions from many brain areas – and that can be parsimoniously explained
499 as being caused and sampled by ‘self-as-agent’. In that case, one might expect high integration in
500 the deepest (highest) parts of the predictive hierarchy; i.e., instantiated in interactions between the
501 default mode network or the salience network (which we note, are also key hubs for interoceptive
502 processing), where the self-model might be instantiated (Margulies et al., 2016). It is also possible
503 that the presence of a higher order integrative component of the larger network is not, in itself,
504 sufficiently integrated to constitute the conscious part of the brain – but that its presence and
505 monitoring allows other parts to be integrated enough to become conscious. The monitoring of the
506 self-model essentially underwrites homeostasis, that is, self-maintenance, which must be tightly
507 related to the exertion of causal power over or the causal constraining of oneself. Indeed, it has
508 been argued by Marshall and colleagues (2017) that intrinsic control and maintenance of causal
509 borders is characteristic of living systems, which seems to align with active inference and
510 predictive coding formulations.

511 Further, IIT’s image of components of a system forming concepts, that in turn can form
512 higher order concepts through integration – for example, integrating the single notes of a song into
513 a melody – might suit such a thing as a self-model particularly well, for the self-model might be
514 thought of as the highest order concept integrating all those lower-order concepts that relate to
515 oneself. This should in particular integrate concepts somehow related to the body, such as
516 interoceptive processes, with perceptual concepts about the current environment as well as those
517 about the agent itself. Finally, it is worth noting that the fact that a higher amount of negative
518 concepts result in higher levels of Φ suits well the argument that counterfactual depth is related to
519 consciousness (Corcoran et al., 2020). Having negative concepts at least conceptually (if not
520 formally) seems related to having a model or experience of the world that describes not only what
521 is, but also what could have been, providing one platform where the otherwise very different
522 theories might meet.

523 Now, how might interoceptive inference fit into IIT’s story of consciousness? Initially, the
524 fit seems poor here as well; for IIT is concerned with the consciousness of systems in general, and
525 additionally also mainly concerned with the experience of these systems *independently* of the
526 external world around them. Interoceptive inference is mainly defined specifically in relation to
527 the brain making inferences about the body within which it is located. We must therefore first
528 allow our conceptualization of interoception to cover any conscious system’s inferences about any
529 kind of body – be it that of a human, animal, plant or complex machine. In addition, we must
530 assume that the conscious experience of a system under IIT must have some kind of relation
531 (structural, perhaps, rather than representational) to the surrounding environment, including the
532 body. This assumption should be treated with caution however, as bridging phenomenological and
533 more functional accounts in this way is no simple project. Here, we offer a speculative outline of
534 these potential links for further discussion.

535 For example, from the view of IIT, one can define the body, generally, as a part of the
536 environment that situates conscious processing, and that it must both react to and control in order
537 to persist, as well as to navigate the rest of the environment. In this view, homeostasis simply
538 becomes acting on or controlling directly that part of the environment that is always present and
539 that I am tightly coupled to, the body, while allostasis is recast as acting on the rest of the
540 environment through the body. It should then be likely that any complex system has in its concept
541 structure some concepts related to bodily states. These concepts need not be about the body *per*
542 *se*; they can be experienced in any way, as long as they are a result of the system navigating within
543 and controlling its body. This means that emotions, understood as embodied-inference (Barrett,
544 2017; Hesp, Smith, et al., 2019), can certainly act as concepts within the system's concept structure
545 that are not in themselves experienced as part of the body, but rather as part of experience itself.
546 One might also hypothesize that conscious systems with complex bodies that need complex
547 behaviour to control and navigate their environment must also be more integrated, and have a
548 richer concept structure that allows for a diverse variety of emotions – in line with how systems
549 become more conscious if they evolve to navigate in a more complex external environment
550 (Albantakis et al., 2014). In this way, a more complex body could directly afford a richer
551 experience with more options and nuances for emotive concepts and the concurrent higher number
552 of negative concepts: a hypothesis that could in principle be investigated in simulation studies. In
553 particular, it may be that high demands on and capabilities for allostasis require a system to be
554 highly integrated and result in a rich bodily experience, since homeostasis by itself is arguably
555 simple.

556 In IIT, conscious experience occurs when a system is able to constrain its own future in a
557 way irreducible to its component elements. Interoceptive inference, then, is inferences about the
558 survival probabilities of the system itself or at least its nearest and most intimate surroundings, the
559 body. Interoceptive inference is therefore crucial for a system to be able to self-constrain in very
560 complex environments. The brain certainly depends on it in order to survive, which can be seen as
561 a type of self-constraining. Successful interoceptive inference may also allow the brain to be more
562 integrated with the body; in IIT terms, that is, to couple with the body in an interdependent way.
563 Given that parts of the brain are so highly integrated that their Φ levels are higher than that between
564 brain and body, probably it is unrealistic (even if theoretically possible) that the brain and body
565 would be so integrated that they together would form one conscious complex; but the adaptive
566 value of high integration will still be in effect even if only a part of the brain stays conscious.

567 One could also imagine that something like a heart – that is, a rhythmical oscillating state
568 which is strongly connected to the rest of the body and brain – would have great effect on a
569 conscious system's concept structure and experience (Allen et al., 2019). It may thus be intriguing
570 to develop evolutionary simulations such as those of Albantakis and colleagues (2014), but with
571 agents that have minimal bodies and task-relevant rhythmically oscillating states that affect the
572 conscious 'brain', too see if such agents evolve concepts relevantly similar to emotions, indicating
573 a phenomenological experience of a bodily state.

574 Notions of selfhood are also important for some theories of consciousness. What might
575 selfhood be in IIT, and would it be related to interoception? Selfhood in IIT could be thought to
576 be the higher-order concept that integrates all body and self-related concepts (emotions, action
577 possibilities and tendencies, in general all that could be called either homeostasis or allostasis).
578 Because the underlying concepts are integrated into a higher order concept, they are not
579 experienced as separate components, but as an integrated whole, meaning that self is the integration
580 of body-related concepts in the same way that a melody is the integration of the single experienced
581 nodes. This might also suggest a fundamental self-other distinction; for it is possible that the
582 components of the system that underwrite bodily and self-experiences are more integrated with
583 each other than they are with experiences related to the external world. This seems likely given
584 that those components might all be influenced by changes in bodily states and therefore be more
585 co-dependent than changes in the external world. Finally, in another sense, there is a ‘self’ in IIT
586 in the sense that there is a main complex which is conscious, a centre of consciousness that is
587 arguably separate from homeostatic and allostatic processes. Would there be a concept within the
588 concept structure specifically related to this - or is it rather the entire concept structure, that is, the
589 entirety of experience as an integrated whole, that might here be called the self, and from which
590 the sense of ‘mineness’ comes? There might here be an opportunity for an I vs. me distinction; i.e.,
591 a distinction between the self as the subject and object within experience (Gallagher, 2000). The
592 former being the entirety of integrated consciousness, and the latter being those concepts in my
593 concept structure that are integrated to form a general experience of what I am and can do,
594 expressly based in bodily experiences. Speculatively, the former might correspond to a lower level
595 primary consciousness – sometimes called C1, for example as in (Frith, 2019) – that does not have
596 meta-representations but is still essentially experienced phenomenologically, while the latter might
597 be a form of higher-order consciousness (C2 and higher, and underwriting access consciousness).

598 Active Inference, Interoception and Consciousness

599 Active inference is a process theory for how adaptive self-organizing systems come to comply
600 with the normative framework of the Free Energy Principle, and thereby stay in existence (K.
601 Friston et al., 2017). There are several theories of how active inference processes might relate to
602 conscious experience; in the following, we first give a brief introduction to active inference under
603 the free energy principle, and then discuss the existing related consciousness theories. Finally, we
604 consider the potential role of interoception within these approaches and active inference in general.

605 The Free Energy Principle (K. Friston, 2010, 2019) is a normative principle, essentially stating
606 any self-organizing system that maintains a non-equilibrium steady state must, in order to resist
607 random perturbations and maintain itself, act as if it minimized its variational free energy, or
608 maximized the Bayesian model evidence, of its implicit model of the world, given sensory
609 observations. This is often situated in an across-scales blanket-oriented formal ontology where
610 reality is described as a nested hierarchy of Markov Blanket structures, that is, statistical

611 separations of internal states from external states (Clark, 2017; Hesp, Ramstead, et al., 2019;
612 Kirchhoff et al., 2018).

613 Blanket states are separated into active states that affect the external world, and sensory
614 states which affect internal states based on impressions from the external world; maintaining a
615 Markov Blanket entails maintaining a non-equilibrium steady state, which mandates gradient
616 flows on variational free energy⁴. These gradient flows mean that, on average, internal states come
617 to statistically model the external world. Furthermore, active states conform to Hamilton's
618 principle of least action so that, on average, active states minimise the path integral of variational
619 free energy over time. Active inference is then a process theory describing *how*, exactly, self-
620 organizing systems might come to minimize their variational free energy now, and in the future
621 (K. Friston et al., 2013; Sajid et al., 2021). On this view, self-organizing systems appear to simulate
622 the consequences of actions in order to select those actions that lead to the least free energy in the
623 future (i.e., least action), leading to a balance between exploratory, information-seeking behaviour,
624 and exploitative, pragmatic behaviour. Active inference (often modelled using Partially
625 Observable Markov Decision Processes) has been used to describe a variety of phenomena,
626 ranging from stratospheric adaption (Rubin et al., 2020) through cellular organization (Kuchling
627 et al., 2020), interoceptive processes (Allen et al., 2019), and neuronal activity (Isomura et al.,
628 2020; Isomura & Friston, 2018) to psychiatric disorders (R. A. Adams et al., 2013; Benrimoh et
629 al., 2018).

630 Active inference is a formal description of how a system interacts with its environment in
631 order to maintain some desired state, and does not necessarily relate inherently to questions about
632 consciousness. There has, however, been work investigating which types of active inference might
633 underlie conscious experience. Most importantly, it has been argued that consciousness is a result
634 of the generative model implied by the system having temporal and counterfactual depth. That is,
635 that it includes future consequences of actions, and that it includes what would have happened had
636 it acted differently in the past (Corcoran et al., 2020; K. Friston, 2018a). Active inference has also
637 been used to answer the meta-hard problem of consciousness (i.e., why we as researchers are so
638 puzzled by the relation between phenomenal experience and reality). It is also argued that complex
639 agents might come to form mid-level beliefs within their hierarchical models of the world as
640 especially certain, but simultaneously come to realize that these beliefs are irreducibly different
641 from the world. This leads to an inferred chasm between the agent's experiences and the external
642 world, and a seeming irreducible difference between subjective experience and objective reality
643 (Clark et al., 2019). Finally, it is argued that the blanket-oriented ontology described before offers
644 a natural separation between intrinsic information geometries on one side, describing how internal
645 states evolve probabilistically over time, and extrinsic information geometries on the other,
646 describing probabilistic beliefs about external states which are then parametrized by internal states,
647 thus uniting the mind/matter distinction under a monist framework (K. J. Friston et al., 2020).

⁴ A gradient flow is simply a description of states that change in the direction of steepest descent on some function of their current value; here, variational free energy.

648 In addition to this, it is theorized that consciousness is the felt effect that results from
649 explicitly evaluating the expected free energy under different actions, as opposed to automatic or
650 reflexive behaviour (Solms & Friston, 2018). There is also an attempt by Ramstead et al.
651 (Ramstead et al., 2021) to apply generative modelling to understand phenomenology on its own
652 terms, arguing that raw sensory experience can be likened to the observations of an active inference
653 agent, and that the coherent lived experience is then the most likely posterior belief or the best
654 explanation for those raw experiences. Many of these approaches are probably consistent or at
655 least overlap with predictive instantiations of HOT and GNWS theories for consciousness in the
656 brain, for they also emphasize hierarchically structured predictions of the consequences of – and
657 the accuracy of – own beliefs. It should also be noted that there are current attempts at synthesizing
658 consciousness theories like Integrated Information Theory and Global Neuronal Workspace theory
659 under the Free Energy Principle to produce a new Integrated World Modelling Theory of
660 consciousness (Safron, 2019, 2020).

661 In this section we take pains to distinguish between active inference, as the general process
662 of acting adaptively by making predicted or preferred states most likely through action in a free
663 energy minimizing fashion, and predictive coding, which is a process that commits to a specific
664 kind of message passing in the brain and how it might come to effectuate such active inference.
665 The two can indeed be closely related, as often seen in the literature, e.g. in (R. Adams et al., 2015),
666 but for clarity we keep them separate for now. This also allows us to distinguish between brain-
667 specific consciousness theories relevant for interoception, and those more general statements about
668 consciousness in complex systems that relate to active inference in general. We focus on the latter
669 here; the former brain-specific predictive processing-based approaches to consciousness have been
670 considered in the previous section.

671 As with the discussion of interoception and IIT above, one can define the body as an
672 external (to the brain or the conscious system, i.e., outside the Markov Blanket) environment that
673 nonetheless is so closely coupled with the brain that it follows it around everywhere, making the
674 body at once both the most important part of the external environment to monitor and predict on
675 one side, and to control on the other. From here it is not a stretch to claim that there can, indeed,
676 in general, not be any successful active inference without at least a rudimentary kind of
677 interoception, for active inference rests on predicting the consequences of one's actions upon the
678 world (and thereby on one's own sensory observations); since the body realizes this influence of
679 actions on the world, then a failure to properly model and make inferences about the body also
680 leads to a (fatal) failure to affect the world in an autopoietic way.

681 This means that successful self-maintaining systems must always model their bodies, be
682 they humans, plants or machines, and that the structure of the body and the actions it can effectuate,
683 therefore, should be strongly determinant for the types of experiences an organism has. One might
684 also consider that the system-within-the-body might model itself *as* the body, that is, in order to
685 reduce unnecessary complexity of its generative model simply coarse-grain itself and its body into
686 a whole in the self-model. This, of course, is only feasible (i.e., free-energy minimizing) if the
687 body and the controlling system (for example the brain) are so tightly coupled that distinguishing

688 between them has only irrelevant advantages to the agent's resulting behaviour. In addition, given
689 that maintenance of the body is crucial for the controlling system's self-maintenance, the
690 pragmatic value – that is, prior preferences over outcomes – for an active inference agent should
691 largely be defined in terms of – or at least in relation to – bodily states, and therefore interoception,
692 largely in accordance with the idea of homeostatic priors and inference as having a privileged
693 position as a first prior (Allen et al., 2020; Allen & Tsakiris, 2018). This clearly posits
694 interoception, self-maintenance, emotion, interoceptive inference, value and consciousness as
695 tightly interlinked concepts.

696 When approaching consciousness and interoception from the perspective of active
697 inference of self-organizing systems in general, rather than specifically from predictive processing
698 in the brain, one might ask why one should focus particularly on the boundary between the brain
699 on the one side, and the body and external world on the other? Markov Blanket partitions can be
700 constructed on many levels of neuronal organization, from single neurons to brain regions (K. J.
701 Friston et al., 2021; Hipólito et al., 2021), indicating that active inference occurs on all those levels,
702 at each level potentially displaying either of the qualities associated with consciousness in the
703 discussion above. Focusing on the level of the brain as a whole situated within the body and the
704 external world would be the traditional choice in consciousness research. It is also the level on
705 which bodily processes such as respiration and heartbeats are part of the immediate environment,
706 and therefore where interoceptive inference happens as it is traditionally conceptualized. Indeed,
707 brain inferences about the body and the world are more similar to personal experiences, compared,
708 for example, to inferences made by a brain region about other brain regions. It might also be
709 hypothesized that the level of the brain as a whole is indeed the level of description with, for
710 example, the longest temporal and counterfactual depth, making it the most interesting level for
711 the purposes of consciousness research.

712 Active inference based accounts of consciousness might be considered functionalistic,
713 because active inference as a framework is centered around how a system interacts with its
714 environment. This is contrasted by Integrated Information Theory (IIT), for example, which
715 focuses on a system's internal causal structure irrespective of its sensorimotor exchanges with the
716 external world. One might imagine two functionally identical (in terms of their blanket states)
717 systems with different internal causal structures, which should therefore have identical generative
718 models from an active inference perspective, but which would have different conscious
719 experiences according to IIT (as in (Oizumi et al., 2014)). A full discussion of the differences
720 between these two approaches are beyond the scope of this paper, but we note that it is a potentially
721 interesting line of research to clarify the theoretical and formal relations between them, for
722 example investigating whether temporal and counterfactual depth of the implied generative model
723 of a system is related to its level of integration.

724 Active inference formalizations have of course already been brought to bear on the question
725 of interoception, in general: volitional control of respiration can be seen as an active inference
726 process which alters interoceptive models (Boyadzhieva & Kayhan, 2020); interoceptive inference
727 has been related to psychopathologies (Paulus et al., 2019); and it lies at the foundation of theories

728 of interoceptive inference in general (A. K. Seth, 2013b). Another recent line of work also tries to
729 understand interoception as a type of active self-inference modulating the volatility of sensory-
730 motor representations. We relate this to consciousness in the following, penultimate section.

731 Interoceptive Self-Inference: an integrated theory of 732 consciousness?

733
734 Finally, we consider how the emerging framework of interoceptive self-inference (Allen, 2020;
735 Allen et al., 2020; Allen & Tsakiris, 2018) might offer an integrative approach to the empirical
736 and theoretical study of consciousness. The theory of interoceptive self-inference is a
737 computational and theoretical model which aims to explain how bodily and interoceptive processes
738 shape exteroceptive and metacognitive awareness, and *vice versa*. Interoceptive self-inference can
739 be seen as a process theory built in part from the FEP, based on empirical and phenomenological
740 observations (Allen & Friston, 2018; Gallagher & Allen, 2018). In particular, the theory posits
741 three core observations:

- 742
743 I. To persist, agents must learn to navigate a volatile, ever-changing world (Piray & Daw,
744 2020a, 2020b; Pulcu & Browning, 2019).
- 745 II. Visceral, homeostatic rhythms directly influence the volatility of both lower-order sensory-
746 motor representations (Allen et al., 2019; Chow et al., 2020; Livneh et al., 2020), and
747 metacognitive inferences thereof (Allen et al., 2016; Hauser et al., 2017).
- 748 III. Therefore, agents actively infer their own volatility trajectories, in part, by sampling and
749 controlling interoceptive rhythms, resulting in close coupling between top-down expected
750 volatility and the visceral body (Lawson et al., 2021; Petzschner et al., 2017, 2021).

751
752 The theory thus proposes that, when estimating our own future reliability or precision, agents
753 intrinsically sample from and predict their own visceral rhythms. Conversely, on shorter
754 timescales, agents can optimize the ‘signal-to-noise’ ratio of ongoing sensorimotor dynamics
755 through ballistic alterations of those same visceral rhythms (Galvez-Pol et al., 2020; Grund et al.,
756 2021). A simple example here is that of the trained sharpshooter, who modulates their breathing
757 in order to align the timing of a trigger pull with the quiescent period. Interoceptive self-inference
758 is thus the implicit, preconscious or prenoetic process by which the confidence and salience of the
759 sensorium is aligned to the rhythms of the body: we literally self-infer our own precision
760 trajectories, and in doing so, we actively shape them.

761 Clearly this process of self-inference aligns closely with philosophical and empirical work
762 which describes the importance of an intrinsic predictive self-model, which contextualizes and
763 embodies phenomenal consciousness (Hohwy, 2016; Limanowski & Blankenburg, 2013; T.
764 Metzinger, 2007). Here we further argue that the minimal-self, i.e., the pre-reflective nature of

765 perceptual consciousness, is closely tied to the interoceptive body, in virtue of the close coupling
766 of these rhythms with the overall stability, reliability, and predictability of the agent's own
767 trajectory. Although the visceral body is rarely the focus of the perceiving self, the interoceptive
768 self-inference model posits that the overall contents of consciousness, and in particular the
769 idiosyncratic salience maps which differ between persons and contexts, are likely to be shaped by
770 the close coupling between expected volatility, sensory-motor precision, and visceral rhythms.
771 Interoceptive self-inference then predicts that sampling of the interoceptive trajectory can be used
772 to estimate the volatility of external states. Cognitive and perceptual biases (e.g., exteroceptive and
773 metacognitive) may then arise from treating interoceptive noise as exteroceptive, such that
774 experimental modulation of interoceptive noise could shift the cognitive bias, as partially
775 demonstrated in recent investigations of interoception and metacognition (Allen et al., 2016;
776 Legrand et al., 2020, 2021) . In parallel, conscious experience may then entail the prioritisation of
777 environmental stimuli which are pertinent to the body's contingencies; for example by increasing
778 the salience of the smell of food when we are hungry.

779 Is interoceptive self-inference then itself an integrative theory of consciousness? Certainly,
780 in light of the previous discussion, we can find links between PHOTT and PGNWS approaches
781 and self-inference. On the self-inference account, the global workspace itself is cast as a dynamic,
782 prospective self-model, which accumulates evidence from cortical and sub-cortical systems to
783 infer an overall estimate of expected precision. Interoceptive prediction errors are thus cast as a
784 controlling factor in the overall bifurcation, topology, and probability-to-ignition of the global
785 workspace. Speculatively, one could potentially re-describe "ignition" as the process of active
786 inference by which a top-down model is self-inferred, meaning, in which the agent engages
787 neuromodulatory and visceromotor processes to actively reshape or reconfigure the overall
788 landscape or topology of precision, literally bringing the moment-to-moment self into existence.
789 This would imply that "ignition" is itself a process of active self-inference, in which the agent
790 entertains one hypothesis over another regarding the overall shape and functionality of the cortical
791 manifold, maintained through the estimation and control of expected precision.

792 Similarly, there are clear potential links between PHOTT and interoceptive-self inference.
793 Interoceptive self-inference was originally developed as a model explaining how and why visceral
794 signals impinge upon metacognitive judgements in other, non-interoceptive domains (Allen et al.,
795 2016; Hauser et al., 2017). Metacognition is typically modelled using a signal-detection theoretic
796 approach, in which subjective confidence or awareness is assumed to depend upon a higher-order
797 meta-representation of first-order signal versus noise, plus some additional metacognitive noise
798 (for review, see the earlier section on PHOTT). Interoceptive self-inference inverts this picture, to
799 suggest that metacognitive estimation is a process of self-inferring the probable correlation
800 between the sensorium and ongoing visceral fluctuations. As a silly example, consider the
801 metacognitive evaluation of whether one will do well on an exam: the confidence estimate here
802 depends both on a judgement of expertise within the domain, and perhaps on whether the agent
803 has been binge-drinking the night before and will thus be suffering from sickness behaviours

804 during the exam. The projection of self-reliability into the future is closely coupled both to domain-
805 relevant knowledge, and the prediction of self-volatility.

806 Interoceptive self-inference would then align itself somewhere between PHOTT and
807 PGNWS, seeking to explain how and why interoceptive prediction errors and precisions are
808 coupled to the cortical hierarchy to shape both top-down predictions of precision, and to actively
809 infer future self-states through descending visceromotor control. However, we wish to pump the
810 brakes a bit here – PGNWS and PHOTT are both currently under-defined process theories. It
811 remains to be seen whether these or any predictive-processing derived theory of consciousness is
812 empirically productive. That is to say, we believe that the ultimate test of a theory of consciousness
813 should not be whether it neatly ties together different conceptual approaches, but whether it can
814 make clear contrasting predictions regarding the mechanisms underlying consciousness itself. And
815 while interoceptive self-inference does make clear empirical predictions about the linkages
816 between say, learning, metacognition, and interoception, it remains to be seen whether these
817 predictions will be similarly fruitful for consciousness research.

818 Conclusion

819 We have reviewed some contemporary approaches to consciousness research in the burgeoning
820 predictive processing literature, with an aim of discovering how research on interoception can
821 inform these emerging discussions. In particular, we highlight links between explanatory concepts
822 found in approaches such as higher-order thought theory, the global neuronal workspace,
823 integrated information theory, and predictive processing versions of these. While our review is by
824 design speculative, we hope to have provided the reader with an overview that can serve as a
825 roadmap for future research in these domains. Overall, we propose that further refinement of the
826 existing theories with consideration for interoceptive inference will prove stimulating to the field.
827 Working out the shared commitments between these different approaches is certainly a
828 monumental endeavour, but one which we hope will ultimately prove fruitful.

829

830 References

- 831 Adams, R. A., Stephan, K. E., Brown, H. R., Frith, C. D., & Friston, K. J. (2013). The
832 Computational Anatomy of Psychosis. *Frontiers in Psychiatry*, 4.
833 <https://doi.org/10.3389/fpsyt.2013.00047>
- 834 Adams, R., Friston, K., & Bastos, A. (2015). Active Inference, Predictive Coding and Cortical
835 Architecture. *Recent Advances On The Modular Organization Of The Cortex*, 97–121.
836 https://doi.org/10.1007/978-94-017-9900-3_7
- 837 Ainley, V., Apps, M. A. J., Fotopoulou, A., & Tsakiris, M. (2016). ‘Bodily precision’: A predictive
838 coding account of individual differences in interoceptive accuracy. *Philosophical
839 Transactions of the Royal Society B: Biological Sciences*, 371(1708), 20160003.
840 <https://doi.org/10.1098/rstb.2016.0003>
- 841 Ainley, V., Tajadura-Jiménez, A., Fotopoulou, A., & Tsakiris, M. (2012). Looking into myself:
842 Changes in interoceptive sensitivity during mirror self-observation. *Psychophysiology*,
843 49(11), 1672–1676.
- 844 Albantakis, L., Hintze, A., Koch, C., Adami, C., & Tononi, G. (2014). Evolution of Integrated
845 Causal Structures in Animats Exposed to Environments of Increasing Complexity. *PLOS
846 Computational Biology*, 10(12), e1003966. <https://doi.org/10.1371/journal.pcbi.1003966>
- 847 Allen, M. (2020). Unravelling the Neurobiology of Interoceptive Inference. *Trends in Cognitive
848 Sciences*, 24(4), 265–266. <https://doi.org/10.1016/j.tics.2020.02.002>
- 849 Allen, M., Frank, D., Schwarzkopf, D. S., Fardo, F., Winston, J. S., Hauser, T. U., & Rees, G.
850 (2016). Unexpected arousal modulates the influence of sensory noise on confidence. *ELife*,
851 5, e18103. PubMed. <https://doi.org/10.7554/eLife.18103>

- 852 Allen, M., & Friston, K. J. (2018). From cognitivism to autopoiesis: Towards a computational
853 framework for the embodied mind. *Synthese*, 195(6), 2459–2482.
854 <https://doi.org/10.1007/s11229-016-1288-5>
- 855 Allen, M., Legrand, N., Correa, C. M. C., & Fardo, F. (2020). Thinking through prior bodies:
856 Autonomic uncertainty and interoceptive self-inference. *Behavioral and Brain Sciences*,
857 43. <https://doi.org/10.1017/S0140525X19002899>
- 858 Allen, M., Levy, A., Parr, T., & Friston, K. J. (2019). In the Body's Eye: The Computational
859 Anatomy of Interoceptive Inference. *BioRxiv*, 603928. <https://doi.org/10.1101/603928>
- 860 Allen, M., & Tsakiris, M. (2018). The body as first prior: Interoceptive predictive processing and
861 the primacy. In *The Interoceptive Mind: From Homeostasis to Awareness* (Vol. 27).
- 862 Apps, M. A., & Tsakiris, M. (2014). The free-energy self: A predictive coding account of self-
863 recognition. *Neuroscience & Biobehavioral Reviews*, 41, 85–97.
- 864 Azzalini, D., Rebollo, I., & Tallon-Baudry, C. (2019). Visceral Signals Shape Brain Dynamics and
865 Cognition. *Trends in Cognitive Sciences*, 23(6), 488–509.
866 <https://doi.org/10.1016/j.tics.2019.03.007>
- 867 Baars, B. J. (1988). *A Cognitive Theory of Consciousness*. Cambridge University Press.
- 868 Baars, B. J., & Franklin, S. (2007). An architectural model of conscious and unconscious brain
869 functions: Global Workspace Theory and IDA. *Neural Networks: The Official Journal of*
870 *the International Neural Network Society*, 20(9), 955–961.
871 <https://doi.org/10.1016/j.neunet.2007.09.013>
- 872 Babo-Rebelo, M., Richter, C. G., & Tallon-Baudry, C. (2016). Neural Responses to Heartbeats in
873 the Default Network Encode the Self in Spontaneous Thoughts. *Journal of Neuroscience*,
874 36(30), 7829–7840. <https://doi.org/10.1523/JNEUROSCI.0262-16.2016>

- 875 Barrett, L. F. (2017). The theory of constructed emotion: An active inference account of
876 interoception and categorization. *Social Cognitive and Affective Neuroscience*, *12*(1), 1–
877 23. <https://doi.org/10.1093/scan/nsw154>
- 878 Barrett, L. F., Quigley, K. S., & Hamilton, P. (2016). An active inference theory of allostasis and
879 interoception in depression. *Philosophical Transactions of the Royal Society B: Biological*
880 *Sciences*, *371*(1708), 20160011. <https://doi.org/10.1098/rstb.2016.0011>
- 881 Barrett, L. F., & Simmons, W. K. (2015). Interoceptive predictions in the brain. *Nature Reviews*.
882 *Neuroscience*, *16*(7), 419–429. <https://doi.org/10.1038/nrn3950>
- 883 Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P., & Friston, K. J. (2012).
884 Canonical Microcircuits for Predictive Coding. *Neuron*, *76*(4), 695–711.
885 <https://doi.org/10.1016/j.neuron.2012.10.038>
- 886 Benrimoh, D., Parr, T., Vincent, P., Adams, R. A., & Friston, K. (2018). Active Inference and
887 Auditory Hallucinations. *Computational Psychiatry (Cambridge, Mass.)*, *2*, 183–204.
888 https://doi.org/10.1162/cpsy_a_00022
- 889 Block, N. (1995). On a Confusion About a Function of Consciousness. *Brain and Behavioral*
890 *Sciences*, *18*(2), 227–247. <https://doi.org/10.1017/s0140525x00038188>
- 891 Boly, M., Massimini, M., Tsuchiya, N., Postle, B. R., Koch, C., & Tononi, G. (2017). Are the
892 Neural Correlates of Consciousness in the Front or in the Back of the Cerebral Cortex?
893 Clinical and Neuroimaging Evidence. *Journal of Neuroscience*, *37*(40), 9603–9613.
894 <https://doi.org/10.1523/JNEUROSCI.3218-16.2017>
- 895 Brown, R., Lau, H., & LeDoux, J. E. (2019). Understanding the Higher-Order Approach to
896 Consciousness. *Trends in Cognitive Sciences*, *0*(0).
897 <https://doi.org/10.1016/j.tics.2019.06.009>

- 898 Bruineberg, J., & Rietveld, E. (2014). Self-organization, free energy minimization, and optimal
899 grip on a field of affordances. *Frontiers in Human Neuroscience*, 8.
900 <https://doi.org/10.3389/fnhum.2014.00599>
- 901 Carruthers, P. (2007). Higher-order theories of consciousness. *The Blackwell Companion to*
902 *Consciousness*, 10, 9780470751466.
- 903 Carruthers, P., & Gennaro, R. (2020). Higher-Order Theories of Consciousness. In E. N. Zalta
904 (Ed.), *The Stanford Encyclopedia of Philosophy* (Fall 2020). Metaphysics Research Lab,
905 Stanford University. [https://plato.stanford.edu/archives/fall2020/entries/consciousness-](https://plato.stanford.edu/archives/fall2020/entries/consciousness-higher/)
906 [higher/](https://plato.stanford.edu/archives/fall2020/entries/consciousness-higher/)
- 907 Chanes, L., & Barrett, L. F. (2016). Redefining the Role of Limbic Areas in Cortical Processing.
908 *Trends in Cognitive Sciences*, 20(2), 96–106. <https://doi.org/10.1016/j.tics.2015.11.005>
- 909 Chow, B. W., Nuñez, V., Kaplan, L., Granger, A. J., Bistrong, K., Zucker, H. L., Kumar, P.,
910 Sabatini, B. L., & Gu, C. (2020). Caveolae in CNS arterioles mediate neurovascular
911 coupling. *Nature*, 579(7797), 106–110. <https://doi.org/10.1038/s41586-020-2026-1>
- 912 Ciaunica, A., Constant, A., Preissl, H., & Fotopoulou, A. (2021). *The First Prior: From Co-*
913 *Embodiment to Co-Homeostasis in Early Life*. PsyArXiv.
914 <https://doi.org/10.31234/osf.io/twubr>
- 915 Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive
916 science. *Behavioral and Brain Sciences*, 36(3), 181–204.
917 <https://doi.org/10.1017/S0140525X12000477>
- 918 Clark, A. (2015). Radical Predictive Processing. *The Southern Journal of Philosophy*, 53(S1), 3–
919 27. <https://doi.org/10.1111/sjp.12120>

- 920 Clark, A. (2017). How to Knit Your Own Markov Blanket. In T. Metzinger & W. Wiese (Eds.),
921 *Philosophy and Predictive Processing*.
- 922 Clark, A., Friston, K., & Wilkinson, S. (2019). Bayesing Qualia Consciousness as Inference, Not
923 Raw Datum. *Journal of Consciousness Studies*, 26.
- 924 Corcoran, A. W., Pezzulo, G., & Hohwy, J. (2020). *From Allostatic Agents to Counterfactual*
925 *Cognisers: Active Inference, Biological Regulation, and the Origins of Cognition*.
926 <https://www.preprints.org/manuscript/201911.0083/v2>
- 927 Critchley, H. D., Wiens, S., Rotshtein, P., Öhman, A., & Dolan, R. J. (2004). Neural systems
928 supporting interoceptive awareness. *Nature Neuroscience*, 7(2), 189–195.
929 <https://doi.org/10.1038/nn1176>
- 930 Dehaene, S., & Changeux, J.-P. (2011). Experimental and Theoretical Approaches to Conscious
931 Processing. *Neuron*, 70(2), 200–227. <https://doi.org/10.1016/j.neuron.2011.03.018>
- 932 Dehaene, S., Changeux, J.-P., Naccache, L., Sackur, J., & Sergent, C. (2006). Conscious,
933 preconscious, and subliminal processing: A testable taxonomy. *Trends in Cognitive*
934 *Sciences*, 10(5), 204–211. <https://doi.org/10.1016/j.tics.2006.03.007>
- 935 Dehaene, S., Lau, H., & Kouider, S. (2017). What is consciousness, and could machines have it?
936 *Science*, 358(6362), 486–492. <https://doi.org/10.1126/science.aan8871>
- 937 Evrard, H. C. (2019). The Organization of the Primate Insular Cortex. *Frontiers in Neuroanatomy*,
938 13. <https://doi.org/10.3389/fnana.2019.00043>
- 939 Feldman, A. G. (2009). New insights into action–perception coupling. *Experimental Brain*
940 *Research*, 194(1), 39–58.
- 941 Feldman, H., & Friston, K. (2010). Attention, Uncertainty, and Free-Energy. *Frontiers in Human*
942 *Neuroscience*, 4. <https://doi.org/10.3389/fnhum.2010.00215>

- 943 Fleming, S. M. (2020). Awareness as inference in a higher-order state space. *Neuroscience of*
944 *Consciousness*, 2020(niz020). <https://doi.org/10.1093/nc/niz020>
- 945 Fleming, S. M., & Daw, N. D. (2016). Self-evaluation of decision-making: A general Bayesian
946 framework for metacognitive computation. *Psychological Review*, 124(1), 91.
947 <https://doi.org/10.1037/rev0000045>
- 948 Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human*
949 *Neuroscience*, 8. <https://doi.org/10.3389/fnhum.2014.00443>
- 950 Fotopoulou, A., & Tsakiris, M. (2017). Mentalizing homeostasis: The social origins of
951 interoceptive inference. *Neuropsychanalysis*, 19(1), 3–28.
952 <https://doi.org/10.1080/15294145.2017.1294031>
- 953 Friston, K. (2009). The free-energy principle: A rough guide to the brain? *Trends in Cognitive*
954 *Sciences*, 13(7), 293–301. <https://doi.org/10.1016/j.tics.2009.04.005>
- 955 Friston, K. (2010). The free-energy principle: A unified brain theory? *Nature Reviews*
956 *Neuroscience*, 11(2), 127–138. <https://doi.org/10.1038/nrn2787>
- 957 Friston, K. (2011). What Is Optimal about Motor Control? *Neuron*, 72(3), 488–498.
958 <https://doi.org/10.1016/j.neuron.2011.10.018>
- 959 Friston, K. (2013). Life as we know it. *Journal of The Royal Society Interface*, 10(86), 20130475.
960 <https://doi.org/10.1098/rsif.2013.0475>
- 961 Friston, K. (2018a). Am i self-conscious? (Or does self-organization entail self-consciousness?).
962 *Frontiers in Psychology*, 9(APR). <https://doi.org/10.3389/fpsyg.2018.00579>
- 963 Friston, K. (2018b). Does predictive coding have a future? *Nature Neuroscience*, 21(8), 1019–
964 1021. <https://doi.org/10.1038/s41593-018-0200-7>

- 965 Friston, K. (2019). A free energy principle for a particular physics. *ArXiv:1906.10184 [q-Bio]*.
966 <http://arxiv.org/abs/1906.10184>
- 967 Friston, K. J., Fagerholm, E. D., Zarghami, T. S., Parr, T., Hipólito, I., Magrou, L., & Razi, A.
968 (2021). Parcels and particles: Markov blankets in the brain. *Network Neuroscience*, 5(1),
969 211–251. https://doi.org/10.1162/netn_a_00175
- 970 Friston, K. J., Wiese, W., & Hobson, J. A. (2020). Sentience and the Origins of Consciousness:
971 From Cartesian Duality to Markovian Monism. *Entropy*, 22(5), 516.
972 <https://doi.org/10.3390/e22050516>
- 973 Friston, K., Lin, M., Frith, C. D., Pezzulo, G., Hobson, J. A., & Ondobaka, S. (2017). Active
974 Inference, Curiosity and Insight. *Neural Computation*, 29(10), 2633–2683.
975 https://doi.org/10.1162/neco_a_00999
- 976 Friston, K., Schwartenbeck, P., Fitzgerald, T., Moutoussis, M., Behrens, T., & Dolan, R. J. (2013).
977 The anatomy of choice: Active inference and agency. *Frontiers in Human Neuroscience*,
978 7. <https://doi.org/10.3389/fnhum.2013.00598>
- 979 Frith, C. D. (2019). The neural basis of consciousness. *Psychological Medicine*, 1–13.
980 <https://doi.org/10.1017/S0033291719002204>
- 981 Gallagher, S. (2000). Philosophical conceptions of the self: Implications for cognitive science.
982 *Trends in Cognitive Sciences*, 4(1), 14–21. [https://doi.org/10.1016/S1364-6613\(99\)01417-](https://doi.org/10.1016/S1364-6613(99)01417-5)
983 5
- 984 Gallagher, S., & Allen, M. (2018). Active inference, enactivism and the hermeneutics of social
985 cognition. *Synthese*, 195(6), 2627–2648. <https://doi.org/10.1007/s11229-016-1269-8>

- 986 Galvez-Pol, A., McConnell, R., & Kilner, J. M. (2020). Active sampling in visual search is coupled
987 to the cardiac cycle. *Cognition*, *196*, 104149.
988 <https://doi.org/10.1016/j.cognition.2019.104149>
- 989 Graziano, M., Guterstam, A., Bio, B., & Wilterson, A. (2019). Toward a standard model of
990 consciousness: Reconciling the attention schema, global workspace, higher-order thought,
991 and illusionist theories. *Cognitive Neuropsychology*, *37*, 1–18.
992 <https://doi.org/10.1080/02643294.2019.1670630>
- 993 Grund, M., Al, E., Pabst, M., Dabbagh, A., Stephani, T., Nierhaus, T., & Villringer, A. (2021).
994 Respiration, heartbeat, and conscious tactile perception. *BioRxiv*, 2021.03.22.436396.
995 <https://doi.org/10.1101/2021.03.22.436396>
- 996 Hauser, T. U., Allen, M., Purg, N., Moutoussis, M., Rees, G., & Dolan, R. J. (2017). Noradrenaline
997 blockade specifically enhances metacognitive performance. *ELife*, *6*.
998 <https://doi.org/10.7554/eLife.24901>
- 999 Hesp, C., Ramstead, M., Constant, A., Badcock, P., Kirchhoff, M., & Friston, K. (2019). A Multi-
1000 scale View of the Emergent Complexity of Life: A Free-Energy Proposal. In G. Y.
1001 Georgiev, J. M. Smart, C. L. Flores Martinez, & M. E. Price (Eds.), *Evolution,*
1002 *Development and Complexity* (pp. 195–227). Springer International Publishing.
1003 https://doi.org/10.1007/978-3-030-00075-2_7
- 1004 Hesp, C., Smith, R., Parr, T., Allen, M., Friston, K., & Ramstead, M. (2019). *Deeply Felt Affect:*
1005 *The Emergence of Valence in Deep Active Inference.* <https://doi.org/10.31234/osf.io/62pfd>
- 1006 Hipólito, I., Ramstead, M. J. D., Convertino, L., Bhat, A., Friston, K., & Parr, T. (2021). Markov
1007 blankets in the brain. *Neuroscience & Biobehavioral Reviews*, *125*, 88–97.
1008 <https://doi.org/10.1016/j.neubiorev.2021.02.003>

- 1009 Hoel, E. P., Albantakis, L., Marshall, W., & Tononi, G. (2016). Can the macro beat the micro?
1010 Integrated information across spatiotemporal scales. *Neuroscience of Consciousness*,
1011 2016(niw012). <https://doi.org/10.1093/nc/niw012>
- 1012 Hohwy, J. (2013). *The Predictive Mind*. Oxford University Press.
- 1013 Hohwy, J. (2016). The Self-Evidencing Brain. *Noûs*, 50(2), 259–285.
1014 <https://doi.org/10.1111/nous.12062>
- 1015 Hohwy, J., & Seth, A. (2020). Predictive processing as a systematic basis for identifying the neural
1016 correlates of consciousness. *Philosophy and the Mind Sciences*, 1(II).
1017 <https://doi.org/10.33735/phimisci.2020.II.64>
- 1018 Isomura, T., & Friston, K. (2018). In vitro neural networks minimise variational free energy.
1019 *Scientific Reports*, 8(1), 16926. <https://doi.org/10.1038/s41598-018-35221-w>
- 1020 Isomura, T., Shimazaki, H., & Friston, K. (2020). Canonical neural networks perform active
1021 inference. *BioRxiv*, 2020.12.10.420547. <https://doi.org/10.1101/2020.12.10.420547>
- 1022 Kanai, R., Komura, Y., Shipp, S., & Friston, K. (2015). Cerebral hierarchies: Predictive
1023 processing, precision and the pulvinar. *Philosophical Transactions of the Royal Society B:*
1024 *Biological Sciences*, 370(1668), 20140169. <https://doi.org/10.1098/rstb.2014.0169>
- 1025 Kilner, J. M., Friston, K. J., & Frith, C. D. (2007). Predictive coding: An account of the mirror
1026 neuron system. *Cognitive Processing*, 8(3), 159–166. [https://doi.org/10.1007/s10339-007-](https://doi.org/10.1007/s10339-007-0170-2)
1027 [0170-2](https://doi.org/10.1007/s10339-007-0170-2)
- 1028 Kirchhoff, M., Parr, T., Palacios, E., Friston, K., & Kiverstein, J. (2018). The Markov blankets of
1029 life: Autonomy, active inference and the free energy principle. *Journal of The Royal*
1030 *Society Interface*, 15(138), 20170792. <https://doi.org/10.1098/rsif.2017.0792>

- 1031 Kleckner, I. R., Zhang, J., Touroutoglou, A., Chanes, L., Xia, C., Simmons, W. K., Quigley, K. S.,
1032 Dickerson, B. C., & Feldman Barrett, L. (2017). Evidence for a large-scale brain system
1033 supporting allostasis and interoception in humans. *Nature Human Behaviour*, *1*(5), 0069.
1034 <https://doi.org/10.1038/s41562-017-0069>
- 1035 Ko, Y., & Lau, H. (2012). A detection theoretic explanation of blindsight suggests a link between
1036 conscious perception and metacognition. *Philosophical Transactions of the Royal Society*
1037 *B: Biological Sciences*, *367*(1594), 1401–1411. <https://doi.org/10.1098/rstb.2011.0380>
- 1038 Koch, C., Massimini, M., Boly, M., & Tononi, G. (2016). Neural correlates of consciousness:
1039 Progress and problems. *Nature Reviews Neuroscience*, *17*(5), 307–321.
1040 <https://doi.org/10.1038/nrn.2016.22>
- 1041 Koster-Hale, J., & Saxe, R. (2013). Theory of Mind: A Neural Prediction Problem. *Neuron*, *79*(5),
1042 836–848. <https://doi.org/10.1016/j.neuron.2013.08.020>
- 1043 Kuchling, F., Friston, K., Georgiev, G., & Levin, M. (2020). Morphogenesis as Bayesian
1044 inference: A variational approach to pattern formation and control in complex biological
1045 systems. *Physics of Life Reviews*, *33*, 88–108. <https://doi.org/10.1016/j.plrev.2019.06.001>
- 1046 Lamme, V. A. F. (2006). Towards a true neural stance on consciousness. *Trends in Cognitive*
1047 *Sciences*, *10*(11), 494–501. <https://doi.org/10.1016/j.tics.2006.09.001>
- 1048 Lamme, V. A. F., & Roelfsema, P. R. (2000). The distinct modes of vision offered by feedforward
1049 and recurrent processing. *Trends in Neurosciences*, *23*(11), 571–579.
1050 [https://doi.org/10.1016/S0166-2236\(00\)01657-X](https://doi.org/10.1016/S0166-2236(00)01657-X)
- 1051 Lau, H. C. (2007). A higher order Bayesian decision theory of consciousness. In R. Banerjee & B.
1052 K. Chakrabarti (Eds.), *Progress in Brain Research* (Vol. 168, pp. 35–48). Elsevier.
1053 [https://doi.org/10.1016/S0079-6123\(07\)68004-2](https://doi.org/10.1016/S0079-6123(07)68004-2)

- 1054 Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious
1055 awareness. *Trends in Cognitive Sciences*, *15*(8), 365–373.
1056 <https://doi.org/10.1016/j.tics.2011.05.009>
- 1057 Lawson, R. P., Bisby, J., Nord, C. L., Burgess, N., & Rees, G. (2021). The Computational,
1058 Pharmacological, and Physiological Determinants of Sensory Learning under Uncertainty.
1059 *Current Biology*, *31*(1), 163-172.e4. <https://doi.org/10.1016/j.cub.2020.10.043>
- 1060 LeDoux, J. E., & Brown, R. (2017). A higher-order theory of emotional consciousness.
1061 *Proceedings of the National Academy of Sciences*, *114*(10), E2016–E2025.
1062 <https://doi.org/10.1073/pnas.1619316114>
- 1063 Limanowski, J. (2017). (Dis-)Attending to the Body. In T. K. Metzinger & W. Wiese (Eds.),
1064 *Philosophy and Predictive Processing*. MIND Group.
1065 <https://doi.org/10.15502/9783958573192>
- 1066 Limanowski, J., & Blankenburg, F. (2013). Minimal self-models and the free energy principle.
1067 *Frontiers in Human Neuroscience*, *7*. <https://doi.org/10.3389/fnhum.2013.00547>
- 1068 Limanowski, J., & Friston, K. (2018). ‘Seeing the Dark’: Grounding Phenomenal Transparency
1069 and Opacity in Precision Estimation for Active Inference. *Frontiers in Psychology*, *9*.
1070 <https://doi.org/10.3389/fpsyg.2018.00643>
- 1071 Livneh, Y., Sugden, A. U., Madara, J. C., Essner, R. A., Flores, V. I., Sugden, L. A., Resch, J. M.,
1072 Lowell, B. B., & Andermann, M. L. (2020). Estimation of Current and Future Physiological
1073 States in Insular Cortex. *Neuron*. <https://doi.org/10.1016/j.neuron.2019.12.027>
- 1074 Mansell, W. (2011). Control of perception should be operationalized as a fundamental property of
1075 the nervous system. *Topics in Cognitive Science*, *3*(2), 257–261.

- 1076 Margulies, D. S., Ghosh, S. S., Goulas, A., Falkiewicz, M., Huntenburg, J. M., Langs, G., Bezgin,
1077 G., Eickhoff, S. B., Castellanos, F. X., Petrides, M., Jefferies, E., & Smallwood, J. (2016).
1078 Situating the default-mode network along a principal gradient of macroscale cortical
1079 organization. *Proceedings of the National Academy of Sciences*, *113*(44), 12574–12579.
1080 <https://doi.org/10.1073/pnas.1608282113>
- 1081 Marshall, W., Kim, H., Walker, S. I., Tononi, G., & Albantakis, L. (2017). How causal analysis
1082 can reveal autonomy in models of biological systems. *Philosophical Transactions of the*
1083 *Royal Society A: Mathematical, Physical and Engineering Sciences*, *375*(2109), 20160358.
1084 <https://doi.org/10.1098/rsta.2016.0358>
- 1085 Mashour, G. A., Roelfsema, P., Changeux, J.-P., & Dehaene, S. (2020). Conscious Processing and
1086 the Global Neuronal Workspace Hypothesis. *Neuron*, *105*(5), 776–798.
1087 <https://doi.org/10.1016/j.neuron.2020.01.026>
- 1088 Metzinger, T. (2007). Empirical perspectives from the self-model theory of subjectivity: A brief
1089 summary with examples. In R. Banerjee & B. K. Chakrabarti (Eds.), *Progress in Brain*
1090 *Research* (Vol. 168, pp. 215–278). Elsevier. [https://doi.org/10.1016/S0079-](https://doi.org/10.1016/S0079-6123(07)68018-2)
1091 [6123\(07\)68018-2](https://doi.org/10.1016/S0079-6123(07)68018-2)
- 1092 Michel, M. (2017). A role for the anterior insular cortex in the global neuronal workspace model
1093 of consciousness. *Consciousness and Cognition*, *49*, 333–346.
1094 <https://doi.org/10.1016/j.concog.2017.02.004>
- 1095 Moran, R. J., Campo, P., Symmonds, M., Stephan, K. E., Dolan, R. J., & Friston, K. J. (2013).
1096 Free Energy, Precision and Learning: The Role of Cholinergic Neuromodulation. *Journal*
1097 *of Neuroscience*, *33*(19), 8227–8236. <https://doi.org/10.1523/JNEUROSCI.4255-12.2013>

- 1098 Odegaard, B., Knight, R. T., & Lau, H. (2017). Should a Few Null Findings Falsify Prefrontal
1099 Theories of Conscious Perception? *Journal of Neuroscience*, 37(40), 9593–9602.
1100 <https://doi.org/10.1523/JNEUROSCI.3217-16.2017>
- 1101 Oizumi, M., Albantakis, L., & Tononi, G. (2014). From the Phenomenology to the Mechanisms
1102 of Consciousness: Integrated Information Theory 3.0. *PLOS Computational Biology*,
1103 10(5), e1003588. <https://doi.org/10.1371/journal.pcbi.1003588>
- 1104 Parr, T., & Friston, K. J. (2017). Working memory, attention, and salience in active inference.
1105 *Scientific Reports*, 7(1), 14678. <https://doi.org/10.1038/s41598-017-15249-0>
- 1106 Parr, T., & Friston, K. J. (2019). Attention or salience? *Current Opinion in Psychology*, 29, 1–5.
- 1107 Paulus, M. P., Feinstein, J. S., & Khalsa, S. S. (2019). An Active Inference Approach to
1108 Interoceptive Psychopathology. *Annual Review of Clinical Psychology*, 15(1), 97–122.
1109 <https://doi.org/10.1146/annurev-clinpsy-050718-095617>
- 1110 Petzschner, F. H., Garfinkel, S. N., Paulus, M. P., Koch, C., & Khalsa, S. S. (2021). Computational
1111 Models of Interoception and Body Regulation. *Trends in Neurosciences*, 44(1), 63–76.
1112 <https://doi.org/10.1016/j.tins.2020.09.012>
- 1113 Petzschner, F. H., Weber, L. A. E., Gard, T., & Stephan, K. E. (2017). Computational
1114 Psychosomatics and Computational Psychiatry: Toward a Joint Framework for Differential
1115 Diagnosis. *Biological Psychiatry*, 82(6), 421–430.
1116 <https://doi.org/10.1016/j.biopsych.2017.05.012>
- 1117 Piray, P., & Daw, N. D. (2020a). A simple model for learning in volatile environments. *PLOS*
1118 *Computational Biology*, 16(7), e1007963. <https://doi.org/10.1371/journal.pcbi.1007963>
- 1119 Piray, P., & Daw, N. D. (2020b). Unpredictability vs. Volatility and the control of learning.
1120 *BioRxiv*, 2020.10.05.327007. <https://doi.org/10.1101/2020.10.05.327007>

- 1121 Pulcu, E., & Browning, M. (2019). The Misestimation of Uncertainty in Affective Disorders.
1122 *Trends in Cognitive Sciences*, 23(10), 865–875. <https://doi.org/10.1016/j.tics.2019.07.007>
- 1123 Quigley, K. S., Kanoski, S., Grill, W. M., Barrett, L. F., & Tsakiris, M. (2021). Functions of
1124 Interoception: From Energy Regulation to Experience of the Self. *Trends in Neurosciences*,
1125 44(1), 29–38. <https://doi.org/10.1016/j.tins.2020.09.008>
- 1126 Ramstead, M. J., Hesp, C., Sandved-Smith, L., Mago, J., Lifshitz, M., Pagnoni, G., Smith, R.,
1127 Dumas, G., Lutz, A., Friston, K., & Constant, A. (2021). *From generative models to*
1128 *generative passages: A computational approach to (neuro)phenomenology*. PsyArXiv.
1129 <https://doi.org/10.31234/osf.io/k9pbn>
- 1130 Rao, R. P. N., & Ballard, D. H. (1999). Predictive coding in the visual cortex: A functional
1131 interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1),
1132 79–87. <https://doi.org/10.1038/4580>
- 1133 Rosenthal, D. M. (2006). Consciousness and Higher-Order Thought. In *Encyclopedia of Cognitive*
1134 *Science*. American Cancer Society. <https://doi.org/10.1002/0470018860.s00149>
- 1135 Rubin, S., Parr, T., Da Costa, L., & Friston, K. (2020). Future climates: Markov blankets and active
1136 inference in the biosphere. *Journal of The Royal Society Interface*, 17(172), 20200503.
1137 <https://doi.org/10.1098/rsif.2020.0503>
- 1138 Safron, A. (2019). *Integrated World Modeling Theory (IWMT) Revisited*. PsyArXiv.
1139 <https://doi.org/10.31234/osf.io/kjngh>
- 1140 Safron, A. (2020). An Integrated World Modeling Theory (IWMT) of Consciousness: Combining
1141 Integrated Information and Global Neuronal Workspace Theories With the Free Energy
1142 Principle and Active Inference Framework; Toward Solving the Hard Problem and

- 1143 Characterizing Agentic Causation. *Frontiers in Artificial Intelligence*, 3.
1144 <https://doi.org/10.3389/frai.2020.00030>
- 1145 Sajid, N., Ball, P. J., Parr, T., & Friston, K. J. (2021). Active inference: Demystified and compared.
1146 *Neural Computation*, 33(3), 674–712. https://doi.org/10.1162/neco_a_01357
- 1147 Seth, A., & Critchley, H. (2013). Extending predictive processing to the body: Emotion as
1148 interoceptive inference. *The Behavioral and Brain Sciences*.
1149 <https://doi.org/10.1017/S0140525X12002270>
- 1150 Seth, A. K. (2013a). Interoceptive inference, emotion, and the embodied self. *Trends in Cognitive*
1151 *Sciences*, 17(11), 565–573. <https://doi.org/10.1016/j.tics.2013.09.007>
- 1152 Seth, A. K. (2013b). Interoceptive inference, emotion, and the embodied self. *Trends in Cognitive*
1153 *Sciences*, 17(11), 565–573. <https://doi.org/10.1016/j.tics.2013.09.007>
- 1154 Seth, A. K., & Friston, K. J. (2016). Active interoceptive inference and the emotional brain.
1155 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371(1708),
1156 20160007. <https://doi.org/10.1098/rstb.2016.0007>
- 1157 Seth, A. K., & Hohwy, J. (2020). Predictive processing as an empirical theory for consciousness
1158 science. *Cognitive Neuroscience*, 0(0), 1–2.
1159 <https://doi.org/10.1080/17588928.2020.1838467>
- 1160 Seth, A. K., & Tsakiris, M. (2018). Being a Beast Machine: The Somatic Basis of Selfhood. *Trends*
1161 *in Cognitive Sciences*, 22(11), 969–981. <https://doi.org/10.1016/j.tics.2018.08.008>
- 1162 Sherrington, C. (1952). *The integrative action of the nervous system*. CUP Archive.
- 1163 Solms, M., & Friston, K. (2018). How and why consciousness arises: Some considerations from
1164 physics and physiology. *Journal of Consciousness Studies*, 25(5–6), 202–238.

- 1165 Sperry, R. W. (1950). Neural basis of the spontaneous optokinetic response produced by visual
1166 inversion. *Journal of Comparative and Physiological Psychology*, 43(6), 482.
- 1167 Sterling, P., & Eyer, J. (1988). Allostasis: A new paradigm to explain arousal pathology. In
1168 *Handbook of life stress, cognition and health* (pp. 629–649). John Wiley & Sons.
- 1169 Synofzik, M., Vosgerau, G., & Newen, A. (2008). Beyond the comparator model: A multifactorial
1170 two-step account of agency. *Consciousness and Cognition*, 17(1), 219–239.
1171 <https://doi.org/10.1016/j.concog.2007.03.010>
- 1172 Tamir, D. I., & Thornton, M. A. (2018). Modeling the Predictive Social Mind. *Trends in Cognitive*
1173 *Sciences*, 22(3), 201–212. <https://doi.org/10.1016/j.tics.2017.12.005>
- 1174 Tononi, G., Boly, M., Massimini, M., & Koch, C. (2016). Integrated information theory: From
1175 consciousness to its physical substrate. *Nature Reviews Neuroscience*, 17(7), 450–461.
1176 <https://doi.org/10.1038/nrn.2016.44>
- 1177 Tononi, G., & Koch, C. (2015). Consciousness: Here, there and everywhere? *Philosophical*
1178 *Transactions of the Royal Society B: Biological Sciences*, 370(1668), 20140167.
1179 <https://doi.org/10.1098/rstb.2014.0167>
- 1180 Vaitl, D. (1996). Interoception. *Biological Psychology*, 42(1), 1–27. <https://doi.org/10.1016/0301->
1181 [0511\(95\)05144-9](https://doi.org/10.1016/0301-0511(95)05144-9)
- 1182 von Helmholtz, H. (1925). Helmholtz’s treatise on physiological optics, (Southall JP, transl.). *New*
1183 *York: Optical Society of America*.
- 1184 Whyte, C. J. (2019). Integrating the global neuronal workspace into the framework of predictive
1185 processing: Towards a working hypothesis. *Consciousness and Cognition*, 73, 102763.
1186 <https://doi.org/10.1016/j.concog.2019.102763>

Interoceptive Inference and Consciousness

- 1187 Whyte, C. J., & Smith, R. (2021). The predictive global neuronal workspace: A formal active
1188 inference model of visual consciousness. *Progress in Neurobiology*, *199*, 101918.
1189 <https://doi.org/10.1016/j.pneurobio.2020.101918>
- 1190 Yeung, N., & Summerfield, C. (2012). Metacognition in human decision-making: Confidence and
1191 error monitoring. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
1192 *367*(1594), 1310–1321. <https://doi.org/10.1098/rstb.2011.0416>
- 1193
1194