

What might interoceptive inference reveal about consciousness?

NIIA NIKOLOVA^{a,1,*}, PETER THESTRUP WAADE^{a,1,2}, KARL J. FRISTON³, AND MICAH ALLEN^{1,4,5}

^aEqual contribution

*Corresponding author: niaa@cfm.au.dk

¹Center for Functionally Integrative Neuroscience, Aarhus University, Denmark

²Interacting Minds Center, Aarhus University, Denmark

³Wellcome Centre for Human Neuroimaging, UCL, United Kingdom

⁴Aarhus Institute of Advanced Studies, Denmark

⁵Cambridge Psychiatry, University of Cambridge, United Kingdom

The mainstream science of consciousness offers a few predominate 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.

INTRODUCTION

What is Consciousness?

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

While these strong contrasts delimit the borders of consciousness, you might also consider the phenomenological properties which reveal themselves upon further reflection. Foremost is the unique “mineness” of any conscious experience. In the transition from sleep to wakefulness, there seem to be distinct properties of ownership and agency. Whereas the infinite void of sleep belongs to no one, even before opening my eyes there is a distinct sense in which experience is happening to someone. In phenomenological terms, we can think about this as the minimal, pre-reflexive conditions about which my experiences are uniquely my own [1]. Consciousness then is something which happens to a sentient subject, which is lived through as the embodied point of view of those seemingly ineffable subjective properties.

A sufficient theory of consciousness then, will deal with each of these properties in turn. What distinguishes conscious states from non-conscious ones? How does selfhood and agency influence these properties? Which sorts of mechanisms give rise to both the phenomenological contents of consciousness, and determines which sorts of states become accessible to conscious

thought? How might the body, or emotion, interact with these properties of consciousness?

Answering these questions is no easy task. Certainly, most who study consciousness have heard the joke that there are as many theories of consciousness as there are consciousness theorists. Our goal here is not to provide a comprehensive predictive processing or active inference theory of consciousness, of which there are already a rapidly growing number (for reviews, see [2–5]). Rather, we aim to illustrate how the notion of interoceptive inference and related concepts might inform the theoretical and empirical science of consciousness, by generating alternative process theories that can then be subject to empirical evaluation.

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

and finally we explore the potential of interoceptive inference to integrate with each of the theories, and how it might illuminate future research directions.

What is Interoceptive Inference?

First, however, we must introduce the standard set pieces of predictive processing and interoceptive inference. Predictive processing can be described as a set of theories which aim to understand how expectations – both neural and psychological – shape, constrain, and ultimately define the mind. These theories have deep roots in cybernetics, information processing, and seminal prospective control models emerging from early 1960s motor and activity theory. A key feature of predictive processing is the basic notion that biological information processing occurs primarily via the minimization of (information theoretic) surprise, such that the nervous system can be understood as a hierarchy of top-down predictions and bottom-up prediction errors. Whilst most early theories extrapolated this basic scheme to explain restricted phenomena such as prospective motor control and the sense of agency [6–8], in recent years these approaches have exploded with a myriad of conceptual, computational, and empirical work. An in-depth review of the scope of predictive processing is beyond this current article. For the unfamiliar reader, we here recall the basic principles, but for a more thorough treatment numerous recent reviews exist, both of the general computational and theoretical principles [9–13], and their relationship with notions of embodiment and selfhood [14–16].

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

Early predictive processing theories largely appealed to this motif of prediction error minimization (PEM) to explain phenomena such as visual perception [17], motor control [18], agency [7], or social cognitive meta-representation [19–21]. In contrast, the new “radical predictive processing” wave embraces the unifying nature of the predictive brain in an attempt to explain how all aspects of information processing and behaviour emerge from the integrated hierarchical flow of predictions, prediction errors and their precision [22, 23]. Within this framework then, we can consider both the specific hierarchical processing of interoceptive sensations [24, 25], and the broader implications of embodied, affective inference with respect to our understanding

¹Invert’ here is using the technical (Bayesian) sense and 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.

of consciousness.

Interoception is generally used to refer to the sensation, perception, and metacognition of the visceral cycles which govern an agent’s homeostasis, allostasis, and ultimately its survival [26–28]. This includes, on the ascending side, the sensory information conveying heartbeats, respiration, and the activity of the stomach and gut to the brain – literally, gut feelings. On the descending side, interoception denotes the visceromotor signals and allostatic reflex arcs by which agents maintain their homeostasis in the face of environmental challenges. Interoceptive processes are thus those which enable an agent to monitor and control the bodily states that are necessary to maintain the balance between energy expenditure and consumption.

We can further demarcate interoceptive processes into those which directly subserve *homeostasis*, that is the maintenance of a steady state defined by specific metabolic set-points, and *allostasis*, the proactive control of the body – and environment – to resolve homeostatic needs before they arise [29–31]. For example, biological necessity dictates that body temperature, blood oxygenation, and blood glucose level are all maintained within a restrictive range of values. Any sustained deviation from these values is likely to negatively impact an organism’s survival, whether through the direct inducement of cellular death, or by the slow attrition of metabolic surplus through starvation. If oxygen is too low, or temperature too high, the brain can directly engage adaptive physiological reflexes, maintaining homeostasis by increasing respiratory frequency or decreasing systolic blood pressure.

These simple sensory-motor reflex arcs, illustrated in Figure 1, can be readily understood by appeal to predictive mechanisms not unlike that of a common household thermostat. That is to say, a low-level spinal, thalamic, or brainstem circuit is generally sufficient to encode the set-point as a prior expectation on the heart-rate, respiratory frequency, or blood pressure. As in afferent control theory, this problem reduces to one of increasing or decreasing the descending visceromotor predictions to minimize any sensory prediction error that occurs: c.f., the equilibrium point hypothesis in motor control [32] and related perceptual control theories [33]. One can thus easily envision simple predictive engrams, which monitor visceral inputs and adjust bodily states as needed to maintain the overall integrity of the system. By comparing the re-afferent sensory inputs to the expected change induced by each top-down prediction, the system can meet whatever thermoregulatory, metabolic, or other homeostatic demands are needed, with relatively little need for higher order cognition.

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

Interoceptive inference in the context of allostasis can thus be viewed as operating at a level once (or thrice) removed from

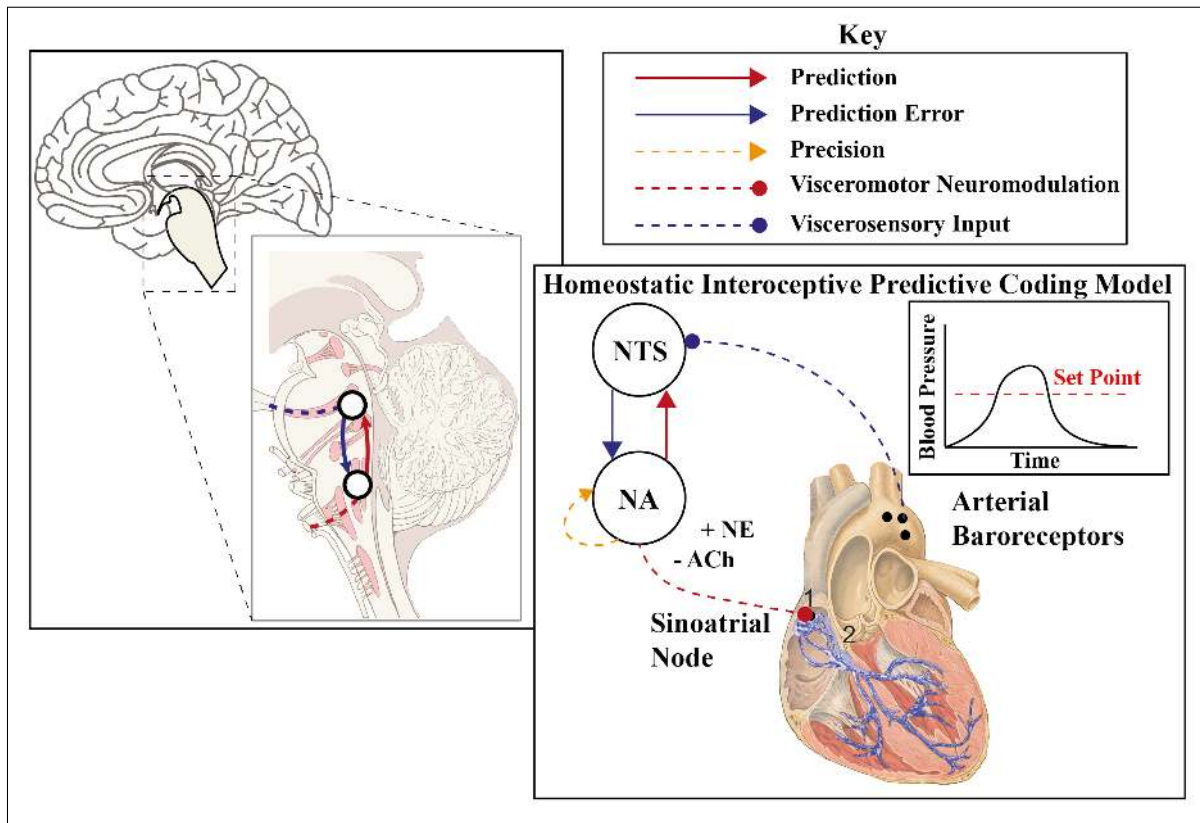


Figure 1. Simplified Homeostatic Control via Interoceptive Inference. This simplified schematic illustrates an example of low-level interoceptive predictive coding in the cardiac domain. Here, a simplified two-node control loop maintains a homeostatic set-point by minimizing the error between afferent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

that of basic homeostasis. Whereas interoceptive inference at the first order might merely involve the regulation of viscerosensory and visceromotor prediction errors, allostatic interoceptive inference requires the agent to link these low-level variables to contextual ones operating at fundamentally longer timescales. As such, interoceptive inference at this level naturally links to the representation of selfhood, valence, and other metacognitive concepts linking the agent's current homeostatic state to the overall volatility of its environment and conspecifics [14, 29, 34].

At a broader level still, we can consider the phylogenetic and ontogenetic role that interoceptive processes play in the overall structure and organization of the nervous system. One standout example of this is found in the Free Energy Principle (FEP), a normative biological theory which posits specific foundational and information theoretic constraints on specific biological process theories [11]. The FEP emphasizes that at the very basis of any biological agent is the self-organized maintenance of its own existence [35]. In this sense, the very structure of the nervous system can be seen as entailing a generative model², which ensures the agent will engage in both homeostatic and allostatic processes. Under the FEP then, the body (both visceral and somato-morphic) are understood as a kind of "first prior" [36, 37], which shapes the evolutionary refinement of the predictive mind. Through this lens, the interoceptive hierarchy plays a special role not only in maintaining an agents' survival, but in determining the salience of every action and ensuing belief updating, and ultimately value itself is understood as whatever maximizes the evidence for the agents' model of a survivable world (c.f., "the self-evidencing brain") [38].

What then can these set pieces about the brain tell us about consciousness? To start, any predictive processing theory will obviously posit a central role for expectations and predictions in the genesis and contents of consciousness. If the mind is primarily concerned with the representation of future events (i.e., the consequences of action), then it seems likely that consciousness is also predominantly prospective. But should a theory of consciousness posit that specific, higher-order modules generate our subjective experience, or rather that it emerges from the collective prediction error minimizing activity of the organism? Similarly, it can be assumed that most predictive processing theories of consciousness should posit a central role in the encoding and modulation of precision – in determining which particular predictions become conscious, and in terms of how conscious predictions should influence affective, metacognitive, and self-related phenomena. That is to say, a basic predictive processing theory of consciousness is likely to ascribe some facets of both *access* and *phenomenal* consciousness (Block, 1995) to error minimizing predictions, and the precision of signals which ensure one particular hypothesis dictates the contents of consciousness versus another.

Do we need appeal to interoceptive processes at all in a theory of consciousness? As we shall see, this depends largely on the overarching theory of consciousness developed, i.e., which conscious phenomena are the target of explanation. Certainly a general PEM-based theory of consciousness would ascribe our bodily self-consciousness to the hierarchical minimization of homeostatic and allostatic prediction errors [39]. E.g., my consciousness of my heart rate or respiration could be argued to

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*.

be a product of the viscerosensory and visceromotor prediction errors and precision signals which drive Bayesian belief updating. In this sense, such interoceptive sensations are likely to dominate my awareness, whenever these systems give off prediction errors, whose precision may need updating³. Yet such a theory would not posit anything particularly unique about interoceptive inference, casting it as just another parcel of the hierarchical organism which gives rise to various bodily aspects of consciousness. Alternatively, one could develop an FEP or similarly radical predictive-processing theory of consciousness, wherein interoceptive inference may fundamentally underpin access and/or phenomenal consciousness. To consider these different possibilities, we now review predominant theories of consciousness in light of interoceptive inference.

PREDICTIVE HIGHER ORDER THOUGHT THEORY (PHOTT)

Higher order-thought theories (HOTT) stem originally from the analytic philosophy of mind [41–43], yet have also found substantive purchase in the empirical science of consciousness [44–46]. In essence, HOTT theories argue that properties of conscious experience arise from the relationship between mental states and higher-order representations of these states [43]. Critically, this implies that a first-order representation by itself is not part of conscious content, unless it is accompanied by another (higher-order) process that is reflecting on its content. In this sense, HOTT stipulates that an agent can be conscious of some representation *X* if and only if the agent possesses a higher order meta-representation of *X*. This approach is based on strong assumptions about the links between phenomenal and access consciousness: according to HOTT, conscious states are by definition those that the agent is aware of.

Empirically speaking, HOTT is often associated with metacognitive approaches to modelling consciousness, such as the popular signal-detection theoretic (SDT) framework [44, 45, 47–49]. Here, for a conscious state to be labelled as such, an experimental subject must not only exhibit above chance accuracy for detecting some stimulus, but also show explicit conscious awareness of their own accuracy, typically measured via subjective confidence or awareness ratings. Now the metric of consciousness is not just whether a subject can reliably discriminate or detect some input, but whether the subject possesses an accurate meta-representation of their own sensory process (i.e., there should also be a high correlation of confidence and accuracy). Neurobiologically, HOTT proponents frequently argue that the prefrontal cortex plays a crucial role in this metacognitive re-representation first-order perceptual contents, and as such is sometimes said to be a necessary and sufficient neural correlate of consciousness (NCC).

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

While thus far no explicit PHOTT theory of consciousness has been proposed, the close alignment of these approaches to empir-

³The updating of the precision of prediction errors is generally read as sensory attention or attenuation [40]. This speaks to an intimate link between conscious (interception-pointing) inference and attentional selection – or sensory attenuation.

ical and computational metacognition research provides some clear starting points. Metacognition, i.e., the meta-representation and control of first order cognitive or perceptual processes, is typically viewed through a decision-theoretic framework in which the agent must monitor the signal and noise distributions underlying first-order perceptual performance, in order to arrive at a representation of the overall probability that one is making correct responses. Fleming ([?]) suggests starting with metacognitive reports of awareness; after all, we can only be aware of another's conscious state through their reports, through language. In the Higher Order State Space (HOSS) model, awareness corresponds to inference on the generative model of the perceptual content, and can be represented as an additional hierarchical state that signals whether perceptual content is present or absent in lower levels. In this case, the higher order thought is cast as a posterior belief over the lower-order contents of consciousness.

Several theorists have proposed Bayesian or predictive variants of these basic models [49–51], where typically the second-order model is seen as integrating the precision of lower-order representations (e.g., the confidence associated with a prediction error encoding a visual input) with high-level “self-priors” describing one's efficacy or overall ability within that cognitive domain. Thus, a basic Bayesian view of metacognition (and meta-representation more generally) posits an extended cognitive hierarchy in which low-level precision signals are read-out and integrated according to some higher-order self-model.

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

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

Alternatively, one could argue for a PHOTT (or perhaps a PHOPT) in which the contents of first-order prediction are largely irrelevant to whether a percept becomes conscious or not, and instead emphasize that PHOTs are fundamentally concerned with meta-representing the precision of lower-order contents. This aligns both with extant Bayesian theories of metacognition, which emphasize that subjective awareness arises from a posterior estimate of precision, and with the intuitive notion that precision is itself fundamentally a second-order statistic (that is, a meta-representation) of first-order predictive processes. In this case, a precision-focused PHOT would likely emphasize the role

of higher-order neural modules in extracting and re-representing the precision (but not the contents) of lower-order predictions, and conscious states would be those associated with the greatest a posteriori precision.

Clearly, these examples are meant to serve as high level outlines illustrating how the set-pieces and explanatory concepts present in predictive processing can be circumscribed within a HOTT of consciousness. Much work remains to be done extrapolating from these basic ideas to a rigorous overall theory. We anticipate that along the way, difficult questions will need to be addressed, concerning for example whether PHOTS are fundamentally concerned with contentful meta-representation, or only with representing the confidence or predictability of first-order processes. One interesting question which emerges immediately, for example, is whether any precision signal could be seen as a sufficient higher-order meta-representation, or whether only higher-order *expected precision* signals would qualify. What we mean is that, according to radical predictive processing theories [23], precision signals can be found at all levels of the central nervous system [36, 52].

At each level of the brain's canonical microcircuitry then, there is a kind of meta-representation encoding the precision of prediction errors arising at that level, and these local precision signals govern the overall flow of contents through the cortical hierarchy. Are these low-level meta-representations sufficient for a content to become conscious? If so, it would appear then that a PHOT theory of consciousness may help to unify recurrent neural processing and HOTT approaches [53, 54], as phenomenal consciousness would emerge from the interaction of local recurrent connections and their associated precision weighting low-level perceptual circuits. In contrast, if it is the explicit representation of expected precision (i.e., top-down, typically polymodal predictions of future changes in lower order precision) that renders a lower state conscious or not, then the resulting PHOT would likely ascribe neuromodulatory circuits and prefrontal modules as fundamental for determining consciousness [55, 56].

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

Alternatively, if the preferred PHOTT emphasizes the role of meta-representations encoding expected precision, then interoceptive processes may play a more constitutive role in determining either phenomenal or access consciousness. Generally speaking, the optimization of expected precision has been proffered as a unifying mechanism by which salience, attention, and high-level self-control emerge [57–61]. Furthermore, the very capacity to supply low levels of hierarchical inference with predictions of precision or predictability has been proposed as a necessary condition for qualitative experience; in the sense of precluding phenomenal transparency [62–64].

This approach views bottom-up and top-down attention as emergent properties of minimizing “precision-prediction errors”,

such that the top-down control of expected precision can selectively enhance or inhibit lower-order percepts. Interoceptive prediction errors and precision thereof are here thought to play a unique role in determining what is salient for an agent in any given context, such that unexpected challenges to homeostasis or allostasis essentially govern the innate value of different outcomes. Computational and conceptual models have expanded on this view to describe a process of metacognitive and interoceptive self-inference, in which the a priori expected precision afforded the homeostatic and allostatic fluctuations is always higher than that of sensory-motor channels [24, 65, 66]. As fluctuations in, for example, blood temperature or arterial pulsation, can directly modulate the noise (i.e., inverse precision) of neuronal circuits in a global fashion [67], then the representation of expected precision is argued to both sample directly from the precision of interoceptive prediction errors, and to utilize descending visceromotor control as a means of optimizing sensory precision.

In PHOTT terms then, this could be taken as an argument that visceral prediction and precision signals play an especially important role in the meta-representation of first-order perceptual contents, such that their subjective salience is largely governed by higher-order thoughts about the interaction between the visceral body and the exteroceptive sensorium. In this sense then, both the “shape” or “contours” of phenomenal consciousness, and the likelihood that a percept becomes conscious (i.e., access consciousness) may depend in part on the top-down meta-representation of expected (interoceptive) precision. Such a process theory would then show some alignment between the PHOTT approach, and recent theoretical proposals suggesting that interoceptive signals may play a fundamental role in shaping the “mineness” or subjective quality of conscious experiences [39, 68–71].

GNWS AND INTEROCEPTIVE INFERENCE

The Global Workspace (GW) theories [72–74] originate from the idea that consciousness arises from processing of information in the brain, and the way in which specific information is selected and broadcast across the brain in order to generate a coherent representation. Here, the brain is composed of a set of specialized, local cortical processing units, which are richly interconnected by excitatory pyramidal neurons spanning between frontal and parietal regions. A piece of information, represented in one or several of the processing units, can cross a threshold and be selected for broadcasting (i.e., amplification) in the process of ‘ignition’, whereby it is simultaneously made available to all processing units. For example, when a bird perches nearby and chirps, my attention is drawn to the sound, and I will gaze around to find the source. The perceptual inputs associated with the bird are carried up and processed, and as they enter the global workspace and become ‘globally available’ as a part of consciousness, such that, they, along with the idea of the bird and the feeling the moment is associated with in my body, are broadcast to various brain systems. These may include memory allowing me to remember the moment, motor action, or higher cognitive systems which enable me to make decisions and talk about my experience. Crucially, most information that is available to and processed by the brain need not enter the global workspace, here consciousness is about how and which information is selected for global processing and awareness. The Global Neuronal Workspace (GNW) theory [74, 75] specifies that information which does become available to the

GNW then recruits brain networks extending over frontal and parietal regions which can integrate the dispersed sources of information into a coherent conscious phenomenon. Thus, the prefrontal cortex plays a central role within GNW as in HOT theories, yet they differ in what functions they ascribe to it [76]; in HOT the higher-order metacognitive processes representing first-order states are what constitute consciousness, so if they are in the PFC, this region becomes a source of consciousness. In GNW meanwhile, conscious states emerge by the broadcasting of information across systems, which can happen due to long range connections between PFC, other fronto-parietal regions comprising the GW. We emphasize that HOTT and GNW are not mutually exclusive, and in fact, several works aim to bridge and unify these theories [77, 78]. The Attentional Schema Theory for example, merges GNW and HOT by proposing that attention amplifies signals so that they may reach ignition, and that there is a higher order representation of the GW which represents the dynamics and implications of having a GW, which is what gives rise to phenomenological consciousness.

Unifying approaches are also building active inference-based models within the GNW framework. The predictive global neuronal workspace (PGNW) [4, 5, 13] combines Bayesian active inference with experimentally corroborated components of the GNW. The PGNW enables us to examine one of the core questions arising from the GNW theories: what determines the ignition threshold? Within predictive processing, the information that crosses the threshold to reach ignition is that which best accommodates PEs throughout the hierarchy, so that the best-fitting (PEM) model of the world is selected and broadcast across systems [79]. Ignition then represents the point at which an evidence accumulation process has reached the threshold where it becomes the most likely explanation of the world (i.e., the current ascending input). The PGNW therefore represents ignition as an inferential process. As in the active inference framework below, ignition here requires sufficient temporal thickness to coordinate and contextualize lower levels of processing [5]. In order to be able to speak of my experience of the chirping bird, I need a representation that is maintained for some period of time and that extends back in time to include me observing the bird. According to the standard GNW account the anterior insula, a key hub processing visceral information and involved in interoceptive awareness [80, 81], selects and prioritizes information prior to possible amplification by the GW [82]. Another theory in the same spirit [83], presents the limbic cortex (including the anterior insula, anterior cingulate cortex, among other areas) as the ‘limbic workspace’ in light of the rich bi-directional connections between these areas and lower levels of processing. In this view, cortical lamination is a distinguishing feature, so that predictions move up from less to more laminated areas while PEs move down in the opposite direction.

Within the PGNW view, interoception is a perceptual system (or set of systems), sensing the internal states and rhythms occurring in the body, and information from it can independently or together with congruent information from other systems, be broadcast by the GNW. For example, I may become aware of a sudden stomach cramp, which incites me to think about what I have eaten earlier in the day. However, recent evidence proposes that interoception might also play a modulating role on other systems [82, 83], whereby interoceptive prediction errors and associated precisions affect the likelihood that other modules are brought into the GW, driving ignition itself through the modulation of salience. It has been suggested that the brain maintains a self-model representing the status of the body, which

is continuously updated to fit ascending interoceptive input by changing interoceptive PEs [26]. Generally, in these accounts, what achieves ignition can be understood as the relationship between the expected (top-down) and sensory (bottom-up) precision, where when the self-model increases the precision of lower order modules, they become better fitting models of the world and are more likely to reach ignition. Further work has proposed that interoception may play a crucial role within the self-model, by either conditioning expected precision [65, 66], or by modulating the degree to which lower-order representations are interpreted as related to the sense of self [25, 62, 69, 84–86].

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

IIT AND INTEROCEPTIVE INFERENCE

The Integrated Information Theory (IIT) [87, 88] of consciousness is an attempt at a formal method for mathematically describing the conscious experience of any given system, agnostic as to all but the causal structure of its substrate. The theory focuses on making an intrinsic description of the system, that is, how the system is to itself, opposed to an extrinsic description from the perspective of an outside observer. The IIT takes as starting point five axioms for what constitutes any phenomenological experience. From that, five criteria are derived which must be met in order for a physical system to support conscious experience: 1) the system must exist, that is, exert and be subject to causal power; and it must do so over itself in a way that is 2) structured of component elements; 3) informative i.e., distinguishable from other causal states; 4) integrated or unitary as a whole, and irreducible to independent subsets; and 5) exclusive or definite, specifying its own borders.

To measure the degree to which a system fulfils these criteria, a measure of *integrated conceptual information* is used, denoted as Φ , which measures the degree to which the system exerts causal power over itself in a way that is irreducible to the activity of its components. The conscious parts of a system are called complexes and are those parts of the system that specify the highest Φ without overlapping with one other. The axiom that complexes cannot overlap also means that smaller complexes are not conscious, even if Φ is larger than zero, as long as they are contained within a larger complex with a higher Φ . Conversely, a large complex is also not conscious if there are smaller complexes within it with a higher Φ . This leads to predictions of consciousness in the brain being situated in areas with more integrated connections, currently thought to be an temporo-parietal-occipital hot-zone in the posterior cortex [89]. This excludes more feed-forward networks like the cerebellum, explaining why this structure does not obviously contribute to consciousness despite its large number of neurons [90]. The exclusivity axiom also means that experience only happens at

one spatial and temporal scale of organization, namely at the level at which Φ is highest [91].

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

The IIT and predictive processing theories of brain function and of consciousness take quite different starting points in a range of respects: IIT is concerned with understanding how systems in general relate to themselves, while predictive processing addresses how the brain, specifically, relates functionally to the surrounding environment, including the body. The former begins entirely in describing phenomenology to identify compatible types of physical systems, while the latter largely takes the opposite direction and starts with what is required for the physical brain in order to self-organize and maintain itself, going from there to describe phenomenology. This makes it challenging to combine the two approaches, a project that is far beyond the scope of this paper – but see [92] for a discussion, in terms of the information geometry of active inference.

It might be worth briefly speculating, however, what predictive processing accounts might be able to offer IIT to inform the broader discussion of interoceptive processing and consciousness. One notion is that hierarchical, precision-weighted prediction error belief updating schemes might provide (neuronal) structures that result in high levels of integration, a suggestion that might potentially be investigated by calculating Φ of canonical neuronal schemas from predictive processing and comparing it to other proposed schemas. The prediction error minimization loops in PEM theories are certainly more complex and integrated than the zero Φ feedforward networks in, say, artificial neural networks. Zooming out, one might also ask if the overall structure of the brain relates to the level of integration; does, for example, the presence of a self-model in the brain somehow constitute or allow for higher levels of integration? One could certainly imagine that the part of the brain that constitutes the ‘self-as-hypothesis’ might be highly integrated, given that it has to coordinate impressions from many brain areas – and that can be parsimoniously explained as being caused and sampled by ‘self-as-agent’. In that case, one might expect high integration in the deepest (highest) parts of the predictive hierarchy; i.e., instantiated in interactions between the default mode network or the salience network (which we note, are also key hubs for interoceptive processing), where the self-model might be instantiated [93]. It is also possible that the presence of a higher order integrative component of the larger network is not, in itself, sufficiently integrated to constitute the conscious part of the brain – but that its presence and monitoring allows other parts to be integrated enough to become conscious. The monitoring of the self-model essentially underwrites homeostasis, that is, self-maintenance, which must be tightly related to the exertion of causal power over or the causal constraining of oneself. Indeed,

it has been argued by Marshall and colleagues [94] that intrinsic control and maintenance of causal borders is characteristic of living systems, which seems to align with active inference and predictive coding formulations.

Further, IIT's image of components of a system forming concepts, that in turn can form higher order concepts through integration – for example, integrating the single notes of a song into a melody – might suit such a thing as a self-model particularly well, for the self-model might be thought of as the highest order concept integrating all those lower-order concepts that relate to oneself. This should in particular integrate concepts somehow related to the body, such as interoceptive processes, with perceptual concepts about the current environment as well as those about the agent itself. Finally, it is worth noting that the fact that a higher amount of negative concepts result in higher levels of Φ suits well the argument that counterfactual depth is related to consciousness [95]. Having negative concepts at least conceptually (if not formally) seems related to having a model or experience of the world that describes not only what is, but also what could have been, providing one platform where the otherwise very different theories might meet.

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

For example, from the view of IIT, one can define the body, generally, as a part of the environment that situates conscious processing, and that it must both react to and control in order to persist, as well as to navigate the rest of the environment. In this view, homeostasis simply becomes acting on or controlling directly that part of the environment that is always present and that I am tightly coupled to, the body, while allostasis is recast as acting on the rest of the environment through the body. It should then be likely that any complex system has in its concept structure some concepts related to bodily states. These concepts need not be about the body *per se*; they can be experienced in any way, as long as they are a result of the system navigating within and controlling its body. This means that emotions, understood as embodied-inference [96, 97], can certainly act as concepts within the system's concept structure that are not in themselves experienced as part of the body, but rather as part of experience itself. One might also hypothesize that conscious systems with complex bodies that need complex behaviour to control and navigate their environment must also be more integrated, and have a richer concept structure that allows for a diverse variety of emotions – in line with how systems become more conscious if they evolve to navigate in a more complex external environment [98]. In this way, a more complex body could directly afford a richer experience with more options and nuances for emotive

concepts and the concurrent higher number of negative concepts: a hypothesis that could in principle be investigated in simulation studies. In particular, it may be that high demands on and capabilities for allostasis require a system to be highly integrated and result in a rich bodily experience, since homeostasis by itself is arguably simple.

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

One could also imagine that something like a heart – that is, a rhythmical oscillating state which is strongly connected to the rest of the body and brain – would have great effect on a conscious system's concept structure and experience [65]. It may thus be intriguing to develop evolutionary simulations such as those of Albantakis and colleagues [98], but with agents that have minimal bodies and task-relevant rhythmically oscillating states that affect the conscious 'brain', to see if such agents evolve concepts relevantly similar to emotions, indicating a phenomenological experience of a bodily state.

Notions of selfhood are also important for some theories of consciousness. What might selfhood be in IIT, and would it be related to interoception? Selfhood in IIT could be thought to be the higher-order concept that integrates all body and self-related concepts (emotions, action possibilities and tendencies, in general all that could be called either homeostasis or allostasis). Because the underlying concepts are integrated into a higher order concept, they are not experienced as separate components, but as an integrated whole, meaning that self is the integration of body-related concepts in the same way that a melody is the integration of the single experienced notes. This might also suggest a fundamental self-other distinction; for it is possible that the components of the system that underwrite bodily and self-experiences are more integrated with each other than they are with experiences related to the external world. This seems likely given that those components might all be influenced by changes in bodily states and therefore be more co-dependent than changes in the external world. Finally, in another sense, there is a 'self' in IIT in the sense that there is a main complex which is conscious, a centre of consciousness that is arguably separate from homeostatic and allostatic processes. Would there be a concept within the concept structure specifically related to this - or is it rather the entire concept structure, that is, the entirety of experience as an integrated whole, that might here be called the self, and from which the sense of 'mineness' comes? There might here be an opportunity for an I vs. me distinction; i.e., a distinction between the self as the subject and object within experience [1]. The former being the entirety of integrated consciousness, and the latter being those concepts in my concept structure that are integrated to form a general experience of what

I am and can do, expressly based in bodily experiences. Speculatively, the former might correspond to a lower level primary consciousness – sometimes called C1, for example as in [99] – that does not have meta-representations but is still essentially experienced phenomenologically, while the latter might be a form of higher-order consciousness (C2 and higher, and underwriting access consciousness).

ACTIVE INFERENCE, INTEROCEPTION AND CONSCIOUSNESS

Active inference is a process theory for how adaptive self-organizing systems come to comply with the normative framework of the Free Energy Principle, and thereby stay in existence [100]. There are several theories of how active inference processes might relate to conscious experience; in the following, we first give a brief introduction to active inference under the free energy principle, and then discuss the existing related consciousness theories. Finally, we consider the potential role of interoception within these approaches and active inference in general.

The Free Energy Principle [101, 102] is a normative principle, essentially stating any self-organizing system that maintains a non-equilibrium steady state must, in order to resist random perturbations and maintain itself, act as if it minimized its variational free energy, or maximized the Bayesian model evidence, of its implicit model of the world, given sensory observations. This is often situated in an across-scales blanket-oriented formal ontology where reality is described as a nested hierarchy of Markov Blanket structures, that is, statistical separations of internal states from external states [103–105].

Blanket states are separated into active states that affect the external world, and sensory states which affect internal states based on impressions from the external world; maintaining a Markov Blanket entails maintaining a non-equilibrium steady state, which mandates gradient flows on variational free energy⁴. These gradient flows mean that, on average, internal states come to statistically model the external world. Furthermore, active states conform to Hamilton's principle of least action so that, on average, active states minimise the path integral of variational free energy over time. Active inference is then a process theory describing *how*, exactly, self-organizing systems might come to minimize their variational free energy now, and in the future [35, 106]. On this view, self-organizing systems appear to simulate the consequences of actions in order to select those actions that lead to the least free energy in the future (i.e., least action), leading to a balance between exploratory, information-seeking behaviour, and exploitative, pragmatic behaviour. Active inference (often modelled using Partially Observable Markov Decision Processes) has been used to describe a variety of phenomena, ranging from stratospheric adaption [107] through cellular organization [108], interoceptive processes [65], and neuronal activity [109, 110] to psychiatric disorders [111, 112].

Active inference is a formal description of how a system interacts with its environment in order to maintain some desired state, and does not necessarily relate inherently to questions about consciousness. There has, however, been work investigating which types of active inference might underlie conscious experience. Most importantly, it has been argued that consciousness is a result of the generative model - entailed by the system - processing

temporal and counterfactual depth. In other words, it models the future consequences of actions, including what would have happened had it acted differently in the past [95, 113]. Active inference has also been used to answer the meta-hard problem of consciousness (i.e., why creatures or researchers are so puzzled by the relation between phenomenal experience and reality). It is also argued that some agents might come to form mid-level beliefs within their hierarchical models of the world as especially certain, but simultaneously come to realize that these beliefs are irreducibly different from the world. This leads to an inferred chasm between the agent's experiences and the external world, and a seeming irreducible difference between subjective experience and objective reality [114]. Finally, it is argued that the blanket-oriented ontology described before offers a natural separation between intrinsic information geometries on one side, describing how internal states evolve probabilistically over time, and extrinsic information geometries on the other, describing probabilistic beliefs about external states which are parametrized by internal states, thus uniting the mind/matter distinction under a monist framework [92].

In addition to this, it is theorized that consciousness is the felt affect that results from explicitly evaluating the expected free energy under different actions, as opposed to automatic or reflexive behaviour [115]. There is also an attempt by Ramstead et al. [116] to apply generative modelling to understand phenomenology on its own terms, arguing that raw sensory experience can be likened to the observations of an active inference agent, and that the coherent lived experience is then the most likely posterior belief or the best explanation for those raw experiences. Many of these approaches are probably consistent or at least overlap with predictive instantiations of HOT and GNWS theories for consciousness in the brain, for they also emphasize hierarchically structured predictions of the consequences of – and the accuracy of – own beliefs. It should also be noted that there are current attempts at synthesizing consciousness theories like Integrated Information Theory and Global Neuronal Workspace theory under the Free Energy Principle to produce a new Integrated World Modelling Theory of consciousness [117, 118].

In this section we take pains to distinguish between active inference, as the general process of acting adaptively by making predicted or preferred states most likely through action in a free energy minimizing fashion, and predictive coding, which is a process that commits to a specific kind of message passing in the brain and how it might come to effectuate such active inference. The two can indeed be closely related, as often seen in the literature, e.g. in [119], but for clarity we keep them separate for now. This also allows us to distinguish between brain-specific consciousness theories relevant for interoception, and those more general statements about consciousness in complex systems that relate to active inference in general. We focus on the latter here; the former brain-specific predictive processing-based approaches to consciousness have been considered in the previous section. As with the discussion of interoception and IIT above, one can define the body as an external (to the brain or the conscious system, i.e., outside the Markov Blanket) environment that nonetheless is so closely coupled with the brain that it follows it around everywhere, making the body at once both the most important part of the external environment to monitor and predict on one side, and to control on the other. From here it is not a stretch to claim that there can, indeed, in general, not be any successful active inference without at least a rudimentary kind of interoception, for active inference rests

⁴A 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.

on predicting the consequences of one's actions upon the world (and thereby on one's own sensory observations); since the body realizes this influence of actions on the world, then a failure to properly model and make inferences about the body also leads to a (fatal) failure to affect the world in an autopoietic way.

This means that successful self-maintaining systems must always model their bodies, be they humans, plants or machines, and that the structure of the body and the actions it can effectuate, therefore, should be strongly determinant for the types of experiences an organism has. One might also consider that the system-within-the-body might model itself *as* the body, that is, in order to reduce unnecessary complexity of its generative model simply coarse-grain itself and its body into a whole in the self-model. This, of course, is only feasible (i.e., free-energy minimizing) if the body and the controlling system (for example the brain) are so tightly coupled that distinguishing between them has only irrelevant advantages to the agent's resulting behaviour. In addition, given that maintenance of the body is crucial for the controlling system's self-maintenance, the pragmatic value – that is, prior preferences over outcomes – for an active inference agent should largely be defined in terms of – or at least in relation to – bodily states, and therefore interoception, largely in accordance with the idea of homeostatic priors and inference as having a privileged position as a first prior [36, 66]. This clearly posits interoception, self-maintenance, emotion, interoceptive inference, value and consciousness as tightly interlinked concepts.

When approaching consciousness and interoception from the perspective of active inference in self-organizing systems in general, rather than specifically from predictive processing in the brain, one might ask why one should focus particularly on the boundary between the brain on the one side, and the body and external world on the other? Markov Blanket partitions can be constructed on many levels of neuronal organization, from single neurons to brain regions [120, 121], indicating that active inference occurs on all those levels—each level potentially displaying either of the qualities associated with consciousness in the discussion above. Focusing on the level of the brain as a whole - situated within the body and the external world—would be the traditional choice in consciousness research. It is also the level on which bodily processes such as respiration and heartbeats are part of the proximal environment, and therefore the provenance of interoceptive inference as typically conceptualized. Indeed, brain inferences about the body and the world underwrite personal experiences - in the sense that the experience is underwritten by inferences about the body—compared, for example, to inferences made by a brain region about other brain regions. It might also be hypothesized that the level of the brain as a whole is indeed the level of description with, for example, the longest temporal and counterfactual depth, making it the most plausible candidate for the purposes of consciousness research. Active inference accounts of consciousness might be considered functionalistic because active inference as a framework is centred around how a system exchanges with its environment. This contrasts with Integrated Information Theory (IIT), for example, which focuses on a system's internal causal structure irrespective of its sensorimotor exchanges with the external world. One might imagine two functionally identical (in terms of their blanket states) systems with different internal causal structures, which reflect different parameterisations of the same generative model from an active inference perspective, but which would have different conscious experiences according to IIT; as in [87]. A full discussion of the differences between

these two approaches are beyond the scope of this paper, but we note that it is a potentially interesting line of research to clarify the theoretical and formal relations between them, for example investigating whether temporal and counterfactual depth of the implied generative model is related to its level of integration.

Active inference formalizations have of course already been brought to bear on the question of interoception, in general: volitional control of respiration can be seen as an active inference process which alters interoceptive models [122]; interoceptive inference has been related to psychopathologies [123]; and it lies at the foundation of theories of interoceptive inference in general [25]. Another recent line of work also tries to understand interoception as a type of active self-inference modulating the volatility of sensory-motor representations. We relate this to consciousness in the following, penultimate section.

INTEROCEPTIVE SELF-INFERENCE: AN INTEGRATED THEORY OF CONSCIOUSNESS?

Finally, we consider how the emerging framework of interoceptive self-inference [24, 36, 66] might offer an integrative approach to the empirical and theoretical study of consciousness. The theory of interoceptive self-inference is a computational and theoretical model which aims to explain how bodily and interoceptive processes shape exteroceptive and metacognitive awareness, and *vice versa*. Interoceptive self-inference can be seen as a process theory built in part from the FEP, based on empirical and phenomenological observations [22, 124]. In particular, the theory posits three core observations:

- I To persist, agents must learn to navigate a volatile, ever-changing world [125–127].
- II Visceral, homeostatic rhythms directly influence the volatility of both lower-order sensory-motor representations [65, 67, 128], and metacognitive inferences thereof [129, 130].
- III Therefore, agents actively infer their own volatility trajectories, in part, by sampling and controlling interoceptive rhythms, resulting in close coupling between top-down expected volatility and the visceral body [14, 34, 131].

The theory thus proposes that, when estimating our own future reliability or precision, agents intrinsically sample from and predict their own visceral rhythms. Conversely, on shorter timescales, agents can optimize the 'signal-to-noise' ratio of ongoing sensorimotor dynamics through ballistic alterations of those same visceral rhythms [132, 133]. A simple example here is that of the trained sharpshooter, who modulates their breathing in order to align the timing of a trigger pull with the quiescent period. Interoceptive self-inference is thus the implicit, preconscious or prenoetic process by which the confidence and salience of the sensorium is aligned to the rhythms of the body: we literally self-infer our own precision trajectories, and in doing so, we actively shape them.

Clearly this process of self-inference aligns closely with philosophical and empirical work which describes the importance of an intrinsic predictive self-model, which contextualizes and embodies phenomenal consciousness [38, 62, 134]. Here we further argue that the minimal-self, i.e., the pre-reflective nature of perceptual consciousness, is closely tied to the interoceptive body, in virtue of the close coupling of these rhythms with the overall stability, reliability, and predictability of the agent's own trajectory. Although the visceral body is rarely

the focus of the perceiving self, the interoceptive self-inference model posits that the overall contents of consciousness, and in particular the idiosyncratic salience maps which differ between persons and contexts, are likely to be shaped by the close coupling between expected volatility, sensory-motor precision, and visceral rhythms. Interoceptive self-inference then predicts that sampling of the interoceptive trajectory can be used to estimate the volatility of external states. Cognitive and perceptual biases (e.g., exteroceptive and metacognitive) may then arise from treating interoceptive noise as exteroceptive, such that experimental modulation of interoceptive noise could shift the cognitive bias, as partially demonstrated in recent investigations of interoception and metacognition [129, 135, 136]. In parallel, conscious experience may then entail the prioritisation of environmental stimuli which are pertinent to the body's contingencies; for example by increasing the salience of the smell of food when we are hungry.

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

Similarly, there are clear potential links between PHOTT and interoceptive-self inference. Interoceptive self-inference was originally developed as a model explaining how and why visceral signals impinge upon metacognitive judgements in other, non-interoceptive domains [129, 130]. Metacognition is typically modelled using a signal-detection theoretic approach, in which subjective confidence or awareness is assumed to depend upon a higher-order meta-representation of first-order signal versus noise, plus some additional metacognitive noise (for review, see the earlier section on PHOTT). Interoceptive self-inference inverts this picture, to suggest that metacognitive estimation is a process of self-inferring the probable correlation between the sensorium and ongoing visceral fluctuations. As a silly example, consider the metacognitive evaluation of whether one will do well on an exam: the confidence estimate here depends both on a judgement of expertise within the domain, and perhaps on whether the agent has been binge-drinking the night before and will thus be suffering from sickness behaviours during the exam. The projection of self-reliability into the future is closely coupled both to domain-relevant knowledge, and the prediction of self-volatility.

Interoceptive self-inference would then align itself somewhere between PHOTT and PGNWS, seeking to explain how and why interoceptive prediction errors and precisions are coupled to the cortical hierarchy to shape both top-down predictions of precision, and to actively infer future self-states through de-

scending visceromotor control. However, we wish to pump the brakes a bit here – PGNWS and PHOTT are both currently under-defined process theories. It remains to be seen whether these or any predictive-processing derived theory of consciousness is empirically productive. That is to say, we believe that the ultimate test of a theory of consciousness should not be whether it neatly ties together different conceptual approaches, but whether it can make clear contrasting predictions regarding the mechanisms underlying consciousness itself. And while interoceptive self-inference does make clear empirical predictions about the linkages between say, learning, metacognition, and interoception, it remains to be seen whether these predictions will be similarly fruitful for consciousness research.

CONCLUSION

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

ACKNOWLEDGEMENTS

Acknowledgments. MA and NN are supported by a Lundbeckfonden Fellowship (under Grant [R272-2017-4345]), and the AIAS-COFUND II fellowship programme that is supported by the Marie Skłodowska-Curie actions under the European Union's Horizon 2020 (under Grant [754513]), and the Aarhus University Research Foundation. KJF was funded by a Wellcome Trust Principal Research Fellowship (Ref: 088130/Z/09/Z).

REFERENCES

1. S. Gallagher, "Philosophical conceptions of the self: implications for cognitive science," *Trends Cogn. Sci.* **4**, 14–21 (2000).
2. J. Hohwy and A. Seth, "Predictive processing as a systematic basis for identifying the neural correlates of consciousness," *Philos. Mind Sci.* **1** (2020). Section: Articles.
3. A. K. Seth and J. Hohwy, "Predictive processing as an empirical theory for consciousness science," *Cogn. Neurosci.* **0**, 1–2 (2020). Publisher: Routledge.
4. C. J. Whyte, "Integrating the global neuronal workspace into the framework of predictive processing: Towards a working hypothesis," *Conscious. Cogn.* **73**, 102763 (2019).
5. C. J. Whyte and R. Smith, "The predictive global neuronal workspace: A formal active inference model of visual consciousness," *Prog. Neurobiol.* **199**, 101918 (2021).
6. R. W. Sperry, "Neural basis of the spontaneous optokinetic response produced by visual inversion." *J. comparative physiological psychology* **43**, 482 (1950). Publisher: American Psychological Association.
7. M. Synofzik, G. Vosgerau, and A. Newen, "Beyond the comparator model: A multifactorial two-step account of agency," *Conscious. Cogn.* **17**, 219–239 (2008).
8. H. von Helmholtz, "Helmholtz's treatise on physiological optics, (southall jp, transl.)," *New York: Opt. Soc. Am.* (1925).

9. A. M. Bastos, W. M. Usrey, R. A. Adams, G. R. Mangun, P. Fries, and K. J. Friston, "Canonical microcircuits for predictive coding," *Neuron* **76**, 695–711 (2012).
10. A. Clark, "Whatever next? predictive brains, situated agents, and the future of cognitive science," *Behav. Brain Sci.* **36**, 181–204 (2013).
11. K. Friston, "The free-energy principle: a rough guide to the brain?" *Trends Cogn. Sci.* **13**, 293–301 (2009).
12. K. Friston, "Does predictive coding have a future?" *Nat. Neurosci.* **21**, 1019–1021 (2018). Number: 8 publisher: Nature Publishing Group.
13. J. Hohwy, *The Predictive Mind* (Oxford University Press, 2013).
14. F. H. Petzschner, S. N. Garfinkel, M. P. Paulus, C. Koch, and S. S. Khalsa, "Computational models of interoception and body regulation," *Trends Neurosci.* **44**, 63–76 (2021).
15. A. Seth and H. Critchley, "Extending predictive processing to the body: emotion as interoceptive inference." *The Behav. brain sciences* (2013).
16. A. K. Seth and K. J. Friston, "Active interoceptive inference and the emotional brain," *Philos. Transactions Royal Soc. B: Biol. Sci.* **371**, 20160007 (2016). Publisher: Royal Society.
17. R. P. N. Rao and D. H. Ballard, "Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects," *Nat. Neurosci.* **2**, 79–87 (1999).
18. K. Friston, "What is optimal about motor control?" *Neuron* **72**, 488–498 (2011).
19. J. M. Kilner, K. J. Friston, and C. D. Frith, "Predictive coding: an account of the mirror neuron system," *Cogn. Process.* **8**, 159–166 (2007).
20. J. Koster-Hale and R. Saxe, "Theory of mind: A neural prediction problem," *Neuron* **79**, 836–848 (2013).
21. D. I. Tamir and M. A. Thornton, "Modeling the predictive social mind," *Trends Cogn. Sci.* **22**, 201–212 (2018).
22. M. Allen and K. J. Friston, "From cognitivism to autopoiesis: towards a computational framework for the embodied mind," *Synthese* **195**, 2459–2482 (2018).
23. A. Clark, "Radical predictive processing," *The South. J. Philos.* **53**, 3–27 (2015).
24. M. Allen, "Unravelling the neurobiology of interoceptive inference," *Trends Cogn. Sci.* **24**, 265–266 (2020).
25. A. K. Seth, "Interoceptive inference, emotion, and the embodied self," *Trends Cogn. Sci.* **17**, 565–573 (2013).
26. L. F. Barrett and W. K. Simmons, "Interoceptive predictions in the brain," *Nat. reviews. Neurosci.* **16**, 419–429 (2015). PMID: 26016744 PMID: PMC4731102.
27. C. Sherrington, *The integrative action of the nervous system* (CUP Archive, 1952).
28. D. Vaitl, "Interoception," *Biol. Psychol.* **42**, 1–27 (1996).
29. L. F. Barrett, K. S. Quigley, and P. Hamilton, "An active inference theory of allostasis and interoception in depression," *Philos. Transactions Royal Soc. B: Biol. Sci.* **371**, 20160011 (2016). Publisher: Royal Society.
30. I. R. Kleckner, J. Zhang, A. Touroutoglou, L. Chanes, C. Xia, W. K. Simmons, K. S. Quigley, B. C. Dickerson, and L. Feldman Barrett, "Evidence for a large-scale brain system supporting allostasis and interoception in humans," *Nat. Hum. Behav.* **1**, 0069 (2017).
31. P. Sterling and J. Eyer, "Allostasis: A new paradigm to explain arousal pathology," in *Handbook of life stress, cognition and health*, (John Wiley Sons, Oxford, England, 1988), pp. 629–649.
32. A. G. Feldman, "New insights into action–perception coupling," *Exp. brain research* **194**, 39–58 (2009). Publisher: Springer.
33. W. Mansell, "Control of perception should be operationalized as a fundamental property of the nervous system," *Top. cognitive science* **3**, 257–261 (2011). Publisher: Wiley Online Library.
34. F. H. Petzschner, L. A. E. Weber, T. Gard, and K. E. Stephan, "Computational psychosomatics and computational psychiatry: Toward a joint framework for differential diagnosis," *Biol. Psychiatry* **82**, 421–430 (2017).
35. K. Friston, P. Schwartenbeck, T. Fitzgerald, M. Moutoussis, T. Behrens, and R. J. Dolan, "The anatomy of choice: active inference and agency," *Front. Hum. Neurosci.* **7** (2013). Publisher: Frontiers.
36. M. Allen and M. Tsakiris, "The body as first prior: Interoceptive predictive processing and the primacy," in *The Interoceptive Mind: From Homeostasis to Awareness*, vol. 27 (2018).
37. A. Ciaunica, A. Constant, H. Preissl, and A. Fotopoulou, "The first prior: from co-embodiment to co-homeostasis in early life," *Tech. rep.* (2021). DOI: 10.31234/osf.io/twubr type: article.
38. J. Hohwy, "The self-evidencing brain," *No@CIRCUMFLEX@us* **50**, 259–285 (2016).
39. V. Ainley, M. A. J. Apps, A. Fotopoulou, and M. Tsakiris, "'bodily precision': a predictive coding account of individual differences in interoceptive accuracy," *Philos. Transactions Royal Soc. B: Biol. Sci.* **371**, 20160003 (2016).
40. H. Brown, R. A. Adams, I. Parees, M. Edwards, and K. Friston, "Active inference, sensory attenuation and illusions," *Cogn. Process.* **14**, 411–427 (2013).
41. P. Carruthers, "Higher-order theories of consciousness," *The Blackwell companion to consciousness* **10**, 9780470751466 (2007). Publisher: Wiley Online Library.
42. P. Carruthers and R. Gennaro, "Higher-order theories of consciousness," in *The Stanford Encyclopedia of Philosophy*, E. N. Zalta, ed. (Metaphysics Research Lab, Stanford University, 2020), fall 2020 ed.
43. D. M. Rosenthal, "Consciousness and higher-order thought," in *Encyclopedia of Cognitive Science*, (American Cancer Society, 2006).
44. R. Brown, H. Lau, and J. E. LeDoux, "Understanding the higher-order approach to consciousness," *Trends Cogn. Sci.* **0** (2019).
45. H. Lau and D. Rosenthal, "Empirical support for higher-order theories of conscious awareness," *Trends Cogn. Sci.* **15**, 365–373 (2011). [Online; accessed 2021-03-03].
46. J. E. LeDoux and R. Brown, "A higher-order theory of emotional consciousness," *Proc. Natl. Acad. Sci.* **114**, E2016–E2025 (2017). Publisher: National Academy of Sciences section: PNAS Plus PMID: 28202735.
47. S. M. Fleming and H. C. Lau, "How to measure metacognition," *Front. Hum. Neurosci.* **8** (2014). [Online; accessed 2019-02-07].
48. Y. Ko and H. Lau, "A detection theoretic explanation of blindsight suggests a link between conscious perception and metacognition," *Philos. Transactions Royal Soc. B: Biol. Sci.* **367**, 1401–1411 (2012). Publisher: Royal Society.
49. H. C. Lau, "A higher order bayesian decision theory of consciousness," in *Progress in Brain Research*, vol. 168 of *Models of Brain and Mind* R. Banerjee and B. K. Chakrabarti, eds. (Elsevier, 2007), pp. 35–48. DOI: 10.1016/S0079-6123(07)68004-2.
50. S. M. Fleming and N. D. Daw, "Self-evaluation of decision-making: A general bayesian framework for metacognitive computation." *Psychol. Rev.* **124**, 91 (2016). [Online; accessed 2019-04-29].
51. N. Yeung and C. Summerfield, "Metacognition in human decision-making: confidence and error monitoring," *Philos. Transactions Royal Soc. B: Biol. Sci.* **367**, 1310–1321 (2012). Publisher: Royal Society.
52. J. Bruineberg and E. Rietveld, "Self-organization, free energy minimization, and optimal grip on a field of affordances," *Front. Hum. Neurosci.* **8** (2014). Publisher: Frontiers.
53. V. A. F. Lamme, "Towards a true neural stance on consciousness," *Trends Cogn. Sci.* **10**, 494–501 (2006). [Online; accessed 2021-03-27].
54. V. A. F. Lamme and P. R. Roelfsema, "The distinct modes of vision offered by feedforward and recurrent processing," *Trends Neurosci.* **23**, 571–579 (2000). [Online; accessed 2021-03-27].
55. M. Boly, M. Massimini, N. Tsuchiya, B. R. Postle, C. Koch, and G. Tononi, "Are the neural correlates of consciousness in the front or in the back of the cerebral cortex? clinical and neuroimaging evidence," *J. Neurosci.* **37**, 9603–9613 (2017). Publisher: Society for Neuroscience section: Dual Perspectives PMID: 28978697.
56. B. Odegaard, R. T. Knight, and H. Lau, "Should a few null findings falsify prefrontal theories of conscious perception?" *J. Neurosci.* **37**, 9593–9602 (2017). Publisher: Society for Neuroscience section: Dual Perspectives PMID: 28978696.
57. H. Feldman and K. Friston, "Attention, uncertainty, and free-energy," *Front. Hum. Neurosci.* **4** (2010). [Online; accessed 2019-03-16].
58. R. Kanai, Y. Komura, S. Shipp, and K. Friston, "Cerebral hierarchies: predictive processing, precision and the pulvinar," *Philos. Transac-*

- tions Royal Soc. B: Biol. Sci. **370**, 20140169 (2015). Publisher: Royal Society.
59. R. J. Moran, P. Campo, M. Symmonds, K. E. Stephan, R. J. Dolan, and K. J. Friston, "Free energy, precision and learning: The role of cholinergic neuromodulation," *J. Neurosci.* **33**, 8227–8236 (2013). Publisher: Society for Neuroscience section: Articles PMID: 23658161.
 60. T. Parr and K. J. Friston, "Working memory, attention, and salience in active inference," *Sci. Reports* **7**, 14678 (2017). [Online; accessed 2019-02-07].
 61. T. Parr and K. J. Friston, "Attention or salience?" *Curr. opinion psychology* **29**, 1–5 (2019). Publisher: Elsevier.
 62. J. Limanowski and F. Blankenburg, "Minimal self-models and the free energy principle," *Front. Hum. Neurosci.* **7** (2013). [Online; accessed 2019-09-20].
 63. J. Limanowski, "(dis-)attending to the body," in *Philosophy and Predictive Processing*, T. K. Metzinger and W. Wiese, eds. (MIND Group, Frankfurt am Main, 2017).
 64. J. Limanowski and K. Friston, "'seeing the dark': Grounding phenomenal transparency and opacity in precision estimation for active inference," *Front. Psychol.* **9** (2018). Publisher: Frontiers.
 65. M. Allen, A. Levy, T. Parr, and K. J. Friston, "In the body's eye: The computational anatomy of interoceptive inference," *bioRxiv* p. 603928 (2019).
 66. M. Allen, N. Legrand, C. M. C. Correa, and F. Fardo, "Thinking through prior bodies: autonomic uncertainty and interoceptive self-inference," *Behav. Brain Sci.* **43** (2020). Publisher: Cambridge University Press.
 67. B. W. Chow, V. Nuñez, L. Kaplan, A. J. Granger, K. Bistrong, H. L. Zucker, P. Kumar, B. L. Sabatini, and C. Gu, "Caveolae in cns arterioles mediate neurovascular coupling," *Nature* **579**, 106–110 (2020). Number: 7797 publisher: Nature Publishing Group.
 68. V. Ainley, A. Tajadura-Jiménez, A. Fotopoulou, and M. Tsakiris, "Looking into myself: Changes in interoceptive sensitivity during mirror self-observation," *Psychophysiology* **49**, 1672–1676 (2012). Publisher: Wiley Online Library.
 69. D. Azzalini, I. Rebollo, and C. Tallon-Baudry, "Visceral signals shape brain dynamics and cognition," *Trends Cogn. Sci.* **23**, 488–509 (2019). [Online; accessed 2019-09-23].
 70. A. Fotopoulou and M. Tsakiris, "Mentalizing homeostasis: The social origins of interoceptive inference," *Neuropsychoanalysis* **19**, 3–28 (2017).
 71. A. K. Seth and M. Tsakiris, "Being a beast machine: The somatic basis of selfhood," *Trends Cogn. Sci.* **22**, 969–981 (2018). [Online; accessed 2019-02-07].
 72. B. J. Baars, *A Cognitive Theory of Consciousness* (Cambridge University Press, 1988).
 73. B. J. Baars and S. Franklin, "An architectural model of conscious and unconscious brain functions: Global workspace theory and ida," *Neural Networks: The Off. J. Int. Neural Netw. Soc.* **20**, 955–961 (2007). PMID: 17998071.
 74. S. Dehaene and J.-P. Changeux, "Experimental and theoretical approaches to conscious processing," *Neuron* **70**, 200–227 (2011). [Online; accessed 2021-03-20].
 75. S. Dehaene, J.-P. Changeux, L. Naccache, J. Sackur, and C. Sergent, "Conscious, preconscious, and subliminal processing: a testable taxonomy," *Trends Cogn. Sci.* **10**, 204–211 (2006). PMID: 16603406.
 76. G. A. Mashour, P. Roelfsema, J.-P. Changeux, and S. Dehaene, "Conscious processing and the global neuronal workspace hypothesis," *Neuron* **105**, 776–798 (2020). [Online; accessed 2021-03-03].
 77. S. Dehaene, H. Lau, and S. Kouider, "What is consciousness, and could machines have it?" *Science* **358**, 486–492 (2017). Publisher: American Association for the Advancement of Science section: Review PMID: 29074769.
 78. M. Graziano, A. Guterstam, B. Bio, and A. Wilterson, "Toward a standard model of consciousness: Reconciling the attention schema, global workspace, higher-order thought, and illusionist theories," *Cogn. Neuropsychol.* **37**, 1–18 (2019).
 79. K. Friston, M. Breakspear, and G. Deco, "Perception and self-organized instability," *Front. Comput. Neurosci.* **6** (2012). Publisher: Frontiers.
 80. H. D. Critchley, S. Wiens, P. Rotshtein, A. Öhman, and R. J. Dolan, "Neural systems supporting interoceptive awareness," *Nat. Neurosci.* **7**, 189–195 (2004). [Online; accessed 2019-08-05].
 81. H. C. Evrard, "The organization of the primate insular cortex," *Front. Neuroanat.* **13** (2019). [Online; accessed 2020-02-05].
 82. M. Michel, "A role for the anterior insular cortex in the global neuronal workspace model of consciousness," *Conscious. Cogn.* **49**, 333–346 (2017). [Online; accessed 2021-03-25].
 83. L. Chanes and L. F. Barrett, "Redefining the role of limbic areas in cortical processing," *Trends Cogn. Sci.* **20**, 96–106 (2016). [Online; accessed 2019-08-29].
 84. M. A. Apps and M. Tsakiris, "The free-energy self: a predictive coding account of self-recognition," *Neurosci. Biobehav. Rev.* **41**, 85–97 (2014).
 85. M. Babo-Rebelo, C. G. Richter, and C. Tallon-Baudry, "Neural responses to heartbeats in the default network encode the self in spontaneous thoughts," *J. Neurosci.* **36**, 7829–7840 (2016). PMID: 27466329.
 86. K. S. Quigley, S. Kanoski, W. M. Grill, L. F. Barrett, and M. Tsakiris, "Functions of interoception: From energy regulation to experience of the self," *Trends Neurosci.* **44**, 29–38 (2021). [Online; accessed 2021-03-26].
 87. M. Oizumi, L. Albantakis, and G. Tononi, "From the phenomenology to the mechanisms of consciousness: Integrated information theory 3.0," *PLOS Comput. Biol.* **10**, e1003588 (2014). Publisher: Public Library of Science.
 88. G. Tononi, M. Boly, M. Massimini, and C. Koch, "Integrated information theory: from consciousness to its physical substrate," *Nat. Rev. Neurosci.* **17**, 450–461 (2016). Number: 7 publisher: Nature Publishing Group.
 89. C. Koch, M. Massimini, M. Boly, and G. Tononi, "Neural correlates of consciousness: progress and problems," *Nat. Rev. Neurosci.* **17**, 307–321 (2016). Number: 5 publisher: Nature Publishing Group.
 90. G. Tononi and C. Koch, "Consciousness: here, there and everywhere?" *Philos. Transactions Royal Soc. B: Biol. Sci.* **370**, 20140167 (2015). Publisher: Royal Society.
 91. E. P. Hoel, L. Albantakis, W. Marshall, and G. Tononi, "Can the macro beat the micro? integrated information across spatiotemporal scales," *Neurosci. Conscious.* **2016** (2016). [Online; accessed 2021-03-26].
 92. K. J. Friston, W. Wiese, and J. A. Hobson, "Sentience and the origins of consciousness: From cartesian duality to markovian monism," *Entropy* **22**, 516 (2020).
 93. D. S. Margulies, S. S. Ghosh, A. Goulas, M. Falkiewicz, J. M. Huntenburg, G. Langs, G. Bezgin, S. B. Eickhoff, F. X. Castellanos, M. Petrides, E. Jefferies, and J. Smallwood, "Situating the default-mode network along a principal gradient of macroscale cortical organization," *Proc. Natl. Acad. Sci.* **113**, 12574–12579 (2016). PMID: 27791099.
 94. W. Marshall, H. Kim, S. I. Walker, G. Tononi, and L. Albantakis, "How causal analysis can reveal autonomy in models of biological systems," *Philos. Transactions Royal Soc. A: Math. Phys. Eng. Sci.* **375**, 20160358 (2017). Publisher: Royal Society.
 95. A. W. Corcoran, G. Pezzulo, and J. Hohwy, "From allostatic agents to counterfactual cognisers: Active inference, biological regulation, and the origins of cognition," (2020). Publisher: Preprints.
 96. L. F. Barrett, "The theory of constructed emotion: an active inference account of interoception and categorization," *Soc. Cogn. Affect. Neurosci.* **12**, 1–23 (2017). [Online; accessed 2021-03-26].
 97. C. Hesp, R. Smith, T. Parr, M. Allen, K. Friston, and M. Ramstead, "Deeply felt affect: The emergence of valence in deep active inference," (2019). Publisher: PsyArXiv.
 98. L. Albantakis, A. Hintze, C. Koch, C. Adami, and G. Tononi, "Evolution of integrated causal structures in animats exposed to environments of increasing complexity," *PLOS Comput. Biol.* **10**, e1003966 (2014). Publisher: Public Library of Science.
 99. C. D. Frith, "The neural basis of consciousness," *Psychol. Medicine* pp. 1–13 (2019). PMID: 31481140.
 100. K. Friston, M. Lin, C. D. Frith, G. Pezzulo, J. A. Hobson, and S. Odojaka, "Active inference, curiosity and insight," *Neural Comput.* **29**, 2633–2683 (2017). [Online; accessed 2019-03-19].

101. K. Friston, "The free-energy principle: a unified brain theory?" *Nat. Rev. Neurosci.* **11**, 127–138 (2010). [Online; accessed 2019-09-20].
102. K. Friston, "A free energy principle for a particular physics," arXiv:1906.10184 [q-bio] (2019). ArXiv: 1906.10184.
103. A. Clark, "How to knit your own markov blanket," in *Philosophy and Predictive Processing*, T. Metzinger and W. Wiese, eds. (2017).
104. C. Hesp, M. Ramstead, A. Constant, P. Badcock, M. Kirchhoff, and K. Friston, "A multi-scale view of the emergent complexity of life: A free-energy proposal," (Springer International Publishing, Cham, 2019), Springer Proceedings in Complexity, pp. 195–227.
105. M. Kirchhoff, T. Parr, E. Palacios, K. Friston, and J. Kiverstein, "The markov blankets of life: autonomy, active inference and the free energy principle," *J. The Royal Soc. Interface* **15**, 20170792 (2018). Publisher: Royal Society.
106. N. Sajid, P. J. Ball, T. Parr, and K. J. Friston, "Active inference: demystified and compared," *Neural Comput.* **33**, 674–712 (2021). ArXiv: 1909.10863.
107. S. Rubin, T. Parr, L. Da Costa, and K. Friston, "Future climates: Markov blankets and active inference in the biosphere," *J. The Royal Soc. Interface* **17**, 20200503 (2020). Publisher: Royal Society.
108. F. Kuchling, K. Friston, G. Georgiev, and M. Levin, "Morphogenesis as bayesian inference: A variational approach to pattern formation and control in complex biological systems," *Phys. Life Rev.* **33**, 88–108 (2020). PMID: 31320316.
109. T. Isomura, H. Shimazaki, and K. Friston, "Canonical neural networks perform active inference," bioRxiv p. 2020.12.10.420547 (2020). Publisher: Cold Spring Harbor Laboratory section: New Results.
110. T. Isomura and K. Friston, "In vitro neural networks minimise variational free energy," *Sci. Reports* **8**, 16926 (2018). Number: 1 publisher: Nature Publishing Group.
111. R. A. Adams, K. E. Stephan, H. R. Brown, C. D. Frith, and K. J. Friston, "The computational anatomy of psychosis," *Front. Psychiatry* **4** (2013). Publisher: Frontiers.
112. D. Benrimoh, T. Parr, P. Vincent, R. A. Adams, and K. Friston, "Active inference and auditory hallucinations," *Comput. Psychiatry (Cambridge, Mass.)* **2**, 183–204 (2018). PMID: 30627670 PMID: PMC6317754.
113. K. Friston, "Am i self-conscious? (or does self-organization entail self-consciousness?)," *Front. Psychol.* **9** (2018). PMID: 29740369 PMID: PMC5928749.
114. A. Clark, K. Friston, and S. Wilkinson, "Bayesing qualia consciousness as inference, not raw datum," *J. Conscious. Stud.* **26** (2019).
115. M. Solms and K. Friston, "How and why consciousness arises: Some considerations from physics and physiology," *J. Conscious. Stud.* **25**, 202–238 (2018). Publisher: Imprint Academic.
116. M. J. Ramstead, C. Hesp, L. Sandved-Smith, J. Mago, M. Lifshitz, G. Pagnoni, R. Smith, G. Dumas, A. Lutz, K. Friston, and A. Constant, "From generative models to generative passages: A computational approach to (neuro)phenomenology," *Tech. rep.* (2021). DOI: 10.31234/osf.io/k9pbn type: article.
117. A. Safron, "Integrated world modeling theory (iwmt) revisited," *Tech. rep.* (2019). DOI: 10.31234/osf.io/kjngh type: article.
118. A. Safron, "An integrated world modeling theory (iwmt) of consciousness: Combining integrated information and global neuronal workspace theories with the free energy principle and active inference framework; toward solving the hard problem and characterizing agentic causation," *Front. Artif. Intell.* **3** (2020). Publisher: Frontiers.
119. R. Adams, K. Friston, and A. Bastos, "Active inference, predictive coding and cortical architecture," *Recent Adv. On The Modul. Organ. Of The Cortex* pp. 97–121 (2015).
120. K. J. Friston, E. D. Fagerholm, T. S. Zarghami, T. Parr, I. Hipólito, L. Magrou, and A. Razi, "Parcels and particles: Markov blankets in the brain," *Netw. Neurosci.* **5**, 211–251 (2021).
121. I. Hipólito, M. J. D. Ramstead, L. Convertino, A. Bhat, K. Friston, and T. Parr, "Markov blankets in the brain," *Neurosci. & Biobehav. Rev.* **125**, 88–97 (2021).
122. A. Boyadzhieva and E. Kayhan, "Keeping the breath in mind: Respiration, neural oscillations and the free energy principle," *PsyArXiv* **4** (2020).
123. M. P. Paulus, J. S. Feinstein, and S. S. Khalsa, "An active inference approach to interoceptive psychopathology," *Annu. Rev. Clin. Psychol.* **15**, 97–122 (2019).
124. S. Gallagher and M. Allen, "Active inference, enactivism and the hermeneutics of social cognition," *Synthese* **195**, 2627–2648 (2018). [Online; accessed 2019-02-07].
125. P. Piray and N. D. Daw, "A simple model for learning in volatile environments," *PLOS Comput. Biol.* **16**, e1007963 (2020). Publisher: Public Library of Science.
126. P. Piray and N. D. Daw, "Unpredictability vs. volatility and the control of learning," bioRxiv p. 2020.10.05.327007 (2020). Publisher: Cold Spring Harbor Laboratory section: New Results.
127. E. Pulcu and M. Browning, "The misestimation of uncertainty in affective disorders," *Trends Cogn. Sci.* **23**, 865–875 (2019). [Online; accessed 2021-03-26].
128. Y. Livneh, A. U. Sugden, J. C. Madara, R. A. Essner, V. I. Flores, L. A. Sugden, J. M. Resch, B. B. Lowell, and M. L. Andermann, "Estimation of current and future physiological states in insular cortex," *Neuron* (2020). [Online; accessed 2020-02-01].
129. M. Allen, D. Frank, D. S. Schwarzkopf, F. Fardo, J. S. Winston, T. U. Hauser, and G. Rees, "Unexpected arousal modulates the influence of sensory noise on confidence," *eLife* **5**, e18103 (2016).
130. T. U. Hauser, M. Allen, N. Purg, M. Moutoussis, G. Rees, and R. J. Dolan, "Noradrenaline blockade specifically enhances metacognitive performance," *eLife* **6** (2017). PMID: 28489001 PMID: PMC5425252.
131. R. P. Lawson, J. Bisby, C. L. Nord, N. Burgess, and G. Rees, "The computational, pharmacological, and physiological determinants of sensory learning under uncertainty," *Curr. Biol.* **31**, 163–172.e4 (2021).
132. A. Galvez-Pol, R. McConnell, and J. M. Kilner, "Active sampling in visual search is coupled to the cardiac cycle," *Cognition* **196**, 104149 (2020).
133. M. Grund, E. Al, M. Pabst, A. Dabbagh, T. Stephani, T. Nierhaus, and A. Villringer, "Respiration, heartbeat, and conscious tactile perception," bioRxiv (2021).
134. T. Metzinger, "Empirical perspectives from the self-model theory of subjectivity: a brief summary with examples," in *Progress in Brain Research*, vol. 168 of *Models of Brain and Mind* R. Banerjee and B. K. Chakrabarti, eds. (Elsevier, 2007), pp. 215–278. DOI: 10.1016/S0079-6123(07)68018-2.
135. N. Legrand, S. S. Engen, C. M. C. Correa, N. K. Mathiasen, N. Nikolova, F. Fardo, and M. Allen, "Emotional metacognition: stimulus valence modulates cardiac arousal and metamemory," *Cogn. Emot.* **0**, 1–17 (2020).
136. N. Legrand, N. Nikolova, C. Correa, M. Braendholt, A. Stuckert, N. Kildahl, M. Vejlo, F. Fardo, and M. Allen, "The heart rate discrimination task: a psychophysical method to estimate the accuracy and precision of interoceptive beliefs," bioRxiv (2021). Publisher: Cold Spring Harbor Laboratory.