

Sussex Research

Theories of consciousness

Anil Seth, Tim Bayne

Publication date

10-06-2023

Licence

This work is made available under the **Copyright not evaluated** licence and should only be used in accordance with that licence. For more information on the specific terms, consult the repository record for this item.

Document Version

Accepted version

Citation for this work (American Psychological Association 7th edition)

Seth, A., & Bayne, T. (2022). *Theories of consciousness* (Version 1). University of Sussex.
<https://hdl.handle.net/10779/uos.23488103.v1>

Published in

Nature Reviews Neuroscience

Link to external publisher version

<https://doi.org/10.1038/s41583-022-00587-4>

Copyright and reuse:

This work was downloaded from Sussex Research Open (SRO). This document is made available in line with publisher policy and may differ from the published version. Please cite the published version where possible. Copyright and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners unless otherwise stated. For more information on this work, SRO or to report an issue, you can contact the repository administrators at sro@sussex.ac.uk. Discover more of the University's research at <https://sussex.figshare.com/>

Theories of consciousness

Anil K. Seth^{1,3*} and Tim Bayne^{2,3,4}

¹Department of Informatics and Sackler Centre for Consciousness Science, University of Sussex, Brighton, United Kingdom.

²School of Philosophical, Historical, and International Studies, Monash University, Melbourne, Victoria, Australia.

³Canadian Institute for Advanced Research (CIFAR) Program on Brain, Mind, and Consciousness, Toronto, Ontario, Canada.

⁴Monash Centre for Consciousness and Contemplative Studies, Monash University, Melbourne, Victoria, Australia

*E-mail: a.k.seth@sussex.ac.uk

Abstract | Recent years have seen a blossoming of theories about the biological and physical basis of consciousness. Good theories guide empirical research, allowing us to interpret data, develop new experimental techniques and expand our capacity to manipulate the phenomenon of interest. Indeed, it is only when couched in terms of a theory that empirical discoveries can ultimately deliver a satisfying understanding of a phenomenon. However, in the case of consciousness, it is unclear how current theories relate to each other, or whether they can be empirically distinguished. To clarify this complicated landscape, we review four prominent theoretical approaches to consciousness: higher-order theories, global workspace theories; reentry and predictive processing theories, and the integrated information theory. We describe the key characteristics of each approach by identifying which aspects of consciousness they propose to explain, what their neurobiological commitments are, and what empirical data are adduced in their support. We consider how some prominent empirical debates might distinguish among these theories, and we outline three ways in which theories need to be developed to deliver a mature regimen of theory-testing in the neuroscience of consciousness. There are good reasons to think that the iterative development, testing and comparison of theories of consciousness will lead to a deeper understanding of this most central of mysteries.

Introduction

In the early decades of its resurgence, the scientific study of consciousness focused on the search for the ‘neural correlates of consciousness’ (NCC). Formally, the NCC of a conscious state are the minimal set of neural events jointly sufficient for that state; in practice the search for NCCs has involved seeking the brain states and processes that are most closely related to consciousness¹⁻³. Focusing on searching for the NCC has been useful because the notion of the NCC is relatively ‘theory neutral’, and thus the NCC framework provided a common language and methodology for researchers with different theoretical and even metaphysical commitments. However, the limitations of the NCC framework have become increasingly clear, as revealed for example in the challenges involved in distinguishing ‘true’ NCCs from the neural pre-requisites and consequences of consciousness⁴⁻⁷. In response to these limitations, there has been a steadily increasing focus on the development of theories of consciousness. With a theory of consciousness (ToC) in hand, we would be able to go beyond a NCC-based methodology and move towards models of consciousness that deliver explanatory insight. Indeed, having an empirically validated ToC should be the primary goal of consciousness science^{8,9}.

Whereas the NCC approach prioritizes the search for correlations between brain activity and consciousness, a theoretical approach instead focuses on identifying explanatory links between neural mechanisms and aspects of consciousness¹⁰. That being said, theorists often employ different conceptions of what it would take to secure an explanatory link between neural activity and consciousness. Some assume that a satisfactory ToC should and can close the ‘explanatory gap’ (**Box 1**), and that it will be possible to render the relationship between neural activity and consciousness as transparent as the relationship between water’s chemical structure and its gross behavioural profile¹¹. Others doubt or remain agnostic as to whether the explanatory gap will ever be fully closed but nonetheless hope for a framework that might explain certain aspects of consciousness and, in doing so, reduce or eliminate the sense of mystery surrounding its biophysical basis^{12,13}. Still others argue that explanatory gap intuitions are misleading, and should not be taken seriously by the science of consciousness.^{14,15}

There is now a wide range of candidate ToCs (**Table 1**). Notably, instead of ToCs progressively being ‘ruled out’ as empirical data accumulates, they seem to be proliferating. This proliferation has led to both attempts to integrate existing theories with each other¹⁶, and to the development of ‘adversarial collaborations’, in which proponents of competing theories agree in advance about whether the outcome of a proposed experiment will support or undermine their preferred theory¹⁷. However, there are significant challenges to both theory integration and adversarial collaboration, as we discuss.

In this Review, we consider how a range of ToCs relate to each other and to empirical data, and we identify some promising avenues by which theory development and empirical research can jointly support each other in the search for a satisfying scientific account of conscious experience. Our attention is restricted to theories that

are either themselves expressed in neurobiological terms, or are plausibly taken to entail claims that can be expressed in neurobiological terms. [As we will see, some ‘neurobiological’ theories of consciousness are expressed in the abstract language of functional relations or information theory, and qualify as ‘neurobiological’ only because the abstract features that they appeal to are associated with particular neural mechanisms.] We also consider only neuroscientific theories that are consistent with known physical theory, and we also leave to one side theories that link consciousness directly to quantum mechanical processes (for examples, see REFS^{18,19}).

Preliminaries

One of the main reasons why ToCs ‘talk’ past each other is that they often have different explanatory targets, for they focus on different aspects of consciousness. We therefore begin by considering what a comprehensive ToC should aim to account for, noting that even this issue is contested, with theorists often disagreeing about what kinds of phenomena a theory of consciousness should explain.

The heart of the problem of consciousness is the issue of ‘experience’ or ‘subjective awareness’. Although no non-circular definition of these terms can be provided, the target phenomenon can be illuminated through some intuitive distinctions. There is ‘something it is like’ for an organism to be conscious²⁰, and what it is like to be in one state of consciousness differs from what it is like to be in another state of consciousness. A comprehensive ToC will explain why some organisms or systems are conscious whereas others are not, and it will explain why states of consciousness differ from each other in the ways that they do.

States of consciousness can be grouped into two classes: global states and local states. Global states concern an organism’s overall subjective profile and are associated with changes in arousal and behavioural responsiveness. Familiar global states include wakefulness, dreaming, sedation, the minimally conscious state, (perhaps) the psychedelic state, and so on. These global states are sometimes called ‘levels’ of consciousness, but we prefer the term ‘global states’ because it leaves open the possibility these states cannot be given a complete ordering in terms of a single dimension but instead are best conceptualized as regions within a multi-dimensional space²¹.

Local states — often referred to as ‘conscious contents’ or as states having ‘qualia’ — are characterized by ‘what it is like’ to be in them. The local state associated with having a headache is distinct from the local state associated with smelling coffee, for what it’s like to have a headache differs from what it’s like to smell coffee. Local states can be described at different levels of granularity, from low-level perceptual features (for example, colour), to objects, to complete multimodal perceptual scenes. An important subset of local states underpins the experience of selfhood, which encompasses experiences of mood, emotion, volition, body-ownership, explicit autobiographical memory, and the like^{13,22-24}. Although neurobiological theories tend to focus on local states with sensory and perceptual content, consciousness also includes local states with cognitive and propositional content, such as the thoughts that arise

when solving a crossword puzzle. Importantly, the local states that an agent has at a particular time do not simply occur as independent elements but are instead bound together as components of a single conscious scene that subsumes each of the agent's local states^{25,26}.

A second distinction is between the phenomenal properties of consciousness and its functional properties. The former term refers to the experiential character of consciousness, as is suggested by the phrase 'what it's like'. The functional aspects of consciousness concern the role(s) that mental states play in the cognitive economy of an organism in virtue of being conscious. ('Function' here encompasses both teleological functions — functional roles as shaped by evolution — and dispositional functions — the role a process plays in the operation of a larger system of which it is a part; see REF²⁷.) For example, being conscious of seeing a coffee cup may enable a range of functions such as the ability to behave flexibly with respect to the cup (perhaps to drink from it, or to throw it across the room), to lay down an episodic memory of the event, to provide verbal reports about the experience, and so on. In making this distinction, we are not claiming that phenomenal and functional properties are independent (they are very likely not independent), merely that they provide distinct explanatory targets for ToCs. As we will see, some ToCs focus on the phenomenal features of consciousness, others focus on the functional features of consciousness, and still others attempt to account for both the functional and phenomenal features of consciousness.

A third distinction is between two kinds of questions concerning local states ('contents') that a ToC might attempt to answer. On one hand, one might ask why an agent is in a certain local state (rather than another). On the other hand, one might ask why a particular local state has the experiential character that it has (rather than an experiential character of some other kind). This distinction can be explained with reference to binocular rivalry, in which each eye is presented with a different stimulus (say, a house to the right eye and a face to the left eye), and the subject's visual experience alternates between the left-eye stimulus and the right-eye stimulus²⁸. Take a particular time at which the contents of consciousness involve a house, while the face is not consciously perceived. Here, we can ask why the mental state corresponding to 'house' is conscious (and that of 'face' is unconscious), and we can also ask why visual experiences of a house have the distinctive experiential character that they have rather than, say, the experiential character of seeing a face, hearing a bell or feeling pain. Notably, there may be some contents that cannot be conscious (for example, low-level processing within early sensory or regulatory systems) and others that can only be conscious (for example, globally integrated perceptual scenes). Thus, in addition to explaining why some mental contents are conscious in some contexts but not others, another challenge is to explain why some contents can never be conscious and why others can exist only as conscious.

Rather than address the full range of issues that we have just identified, most ToCs aim to explain only certain aspects of consciousness, perhaps as a step on the way to becoming comprehensive. Although being restricted in some way is not itself an objection to a ToC, it does mean that the task of inter-theory comparison is less

straightforward than it might otherwise be. If theories are targeting different aspects of consciousness (say, one theory is focused on the phenomenal character of consciousness and another is focused on its functional profile) then they might not be the ‘adversaries’ that they at first glance appear to be.

The ToCs we review here are grouped into four categories: higher-order theories (HOTs), global workspace theories (GWTs), integrated information theory (IIT) and reentry and predictive processing (PP) theories. Although some accounts of consciousness straddle multiple categories, and others are not plausibly subsumed under any of these categories (**Table 1**), this four-way distinction between ToCs provides a useful lens through which to view the current state-of-play in the science of consciousness (**Box 2**; for other ways of grouping theories, see for example REF.²⁹). In what follows, we introduce the key elements of each category; describe some notable within-category differences, and identify those aspects of consciousness most closely associated with each category. We then illustrate how these ToCs relate to each other in terms of some prominent empirical debates, and present several proposals that, we suggest, will help drive a virtuous cycle between theory development and experimental investigation.

Higher-order theories

The core claim uniting all HOTs is that a mental state is conscious in virtue of being the target of a certain kind of meta-representational state. Meta-representations are not merely representations that occur higher or deeper in a processing hierarchy but are rather representations that have as their targets other representations (**Fig. 1**). For example, a representation with the content <I have a visual experience of a moving dot> is a meta-representation, for its content concerns the agent’s own representations of the world rather than the world itself.

An important respect in which HOTs differ from each other concerns the account that they give of the nature and role of the meta-representations that are responsible for consciousness. Some versions of the approach identify the kinds of meta-representations that are crucial for consciousness with thoughts (or thought-like states) that have conceptual content^{30,31 32}. Other varieties of HOT have been expressed in computational terms. According to the self-organizing metarepresentational account, consciousness involves higher-order brain networks learning to re-describe the representations encoded in lower-order networks in a way that counts as meta-representational^{33,34}. Alternatively, higher-order state space theory proposes that subjective reports (for example, statements such as ‘I am aware of X’) are metacognitive (higher-order) decisions about a generative model of perceptual content³⁵, while perceptual reality monitoring posits that conscious perception arises when a higher-order network judges a first-order representation to be a reliable reflection of the external world^{36,37}.

As should be clear from the foregoing, HOTs focus on explaining why some contents are conscious whereas others are not. However, these theories are not limited to this particular focus — they also have the resources to address issues pertaining to

the experiential character of local states. One prominent example concerns the (debated) intuition that the contents of perceptual experience often outstrip the information available in ‘first-order’ sensory representations, as is alleged to occur in the context of peripheral vision^{38,39}. The HOT-based proposal here is that the apparently ‘inflated’ phenomenology of peripheral visual experience is caused by the higher-order misrepresentation of first-order states⁴⁰. The HOT approach can also be extended to explain why some contents are unable to be conscious (they cannot be the targets of appropriate meta-representational states) and why some contents are necessarily conscious (they are necessarily accompanied by appropriate meta-representational states). HOTs rarely focus on global states of consciousness, but it would be natural for them to appeal to the integrity of representational processes to account for the distinctions between global states.

A particularly intriguing question is whether (and if so, how) HOTs explain the distinctive phenomenal character of various kinds of experiences. Why is the phenomenal character associated with seeing a sunset so different from the phenomenal character associated with a headache? The general shape of the higher-order response to this question is that the phenomenal character of a state is determined by the properties that the relevant meta-representational state ascribes to it. Most examples of this approach focus on visual experience⁴⁰, but there have also been attempts to account for the phenomenal character of emotional states⁴¹ and metacognitive states, such as the ‘what it’s like’ to feel confident in a perceptual decision^{42,43}. Ultimately, any fully reductive version of the higher-order approach must explain why the representation of various properties generates the phenomenology that it does (or is identical to it), and how neural activity enables the relevant properties to be represented in the first place.

Higher-order accounts of consciousness are primarily accounts of what makes a mental state conscious and, as such, the approach is not committed to any particular view of the function(s) of consciousness. Indeed, some HOTs downplay the idea that consciousness has any distinctive function⁴⁴. Other versions of the higher-order approach identify the functional role of consciousness with the metacognitive processes associated with confidence judgements and error monitoring⁴⁵. However, while higher-order views allow that conscious mental states can be accompanied by conscious metacognitive judgements — such as those involved in explicit performance monitoring or subjective confidence reports — most versions of this approach do not require that conscious perception is always accompanied by a corresponding conscious metacognitive state. Instead, for meta-representations to be conscious, they themselves must be the objects of a suitable meta-representational state.

With respect to the neural basis of consciousness, the emphasis on meta-representation has led higher-order theorists to emphasize anterior cortical regions, especially the prefrontal cortex³⁰, given the association of these regions with complex cognitive functions. However, although most HOTs propose that anterior involvement is implicated in consciousness, there is disagreement about precisely which anterior regions (or processes) are required⁴⁶.

Global workspace theories

GWTs originate from ‘blackboard’ architectures in artificial intelligence, in which the blackboard is a centralised resource through which specialised processors share and receive information. The first GWT of consciousness⁴⁷ was framed at a cognitive level. It proposed that conscious mental states are those which are ‘globally available’ to a wide range of cognitive processes including attention, evaluation, memory and verbal report. The core claim of GWTs is that it is the wide accessibility of information to such consumer cognitive systems that constitutes conscious experience (**Fig. 2**).

This basic claim has since been developed into a neural theory — often referred to as the “global neuronal workspace theory” — according to which sensory information gains access to consciousness when it is ‘broadcast’ within an anatomically widespread neuronal workspace that is implemented across higher-order cortical association areas, with a particular (although not exclusive) emphasis on the prefrontal cortex^{48,49}. Access to the global workspace is achieved through non-linear network ‘ignition’ in which recurrent processing amplifies and sustains neuronal representations⁵⁰. The emphasis on ignition and broadcast — as compared with meta-representation — is one way in which GWT is distinguished from the HOT approach.

Like HOTs, GWTs focus on the question of what makes a representation conscious, and GWT theorists have rarely attempted to account for the phenomenal differences between distinct kinds of experiences (although though see REF.⁵¹). Returning to our example of binocular rivalry, the GWT view aims to explain why, at a particular point in time, the mental state corresponding to ‘house’ is conscious (while that corresponding to ‘face’ is unconscious), but it offers no direct account of the experiential contrast between seeing a house on the one hand and seeing a face on the other.

The relative silence of GWTs on the issue of experiential character aligns with the general tendency of such theories to focus on functional, rather than phenomenal, aspects of consciousness. In fact, GWTs are often explicitly proposed as accounts of ‘conscious access’⁴⁹; that is, as accounts of why certain representations are available to be flexibly used by a wide range of consuming systems (whereas others aren’t). The core functional property addressed by GWTs is the ability of conscious states to guide behaviour and cognition in flexible, context-dependent ways. GWTs also offer clear accounts of how consciousness is related to other cognitive processes, such as attention and working memory. According to GWTs, attention selects and amplifies specific signals, allowing them to enter the workspace (and thus be conscious); while consciousness and working memory are intimately related because attended working memory items are conscious and use the global workspace for broadcast⁴⁹.

GWTs account for changes in global states of consciousness in terms of alterations to the functional integrity of the workspace. Neurally, a global loss of consciousness is reflected in impaired functional or dynamical connectivity in frontoparietal regions that are considered as ‘hub’ nodes in the global workspace⁵², and

in functional connectivity becoming increasingly constrained to patterns directly reflecting the underlying structural connectivity⁵³⁻⁵⁵.

One important question raised by GWTs concerns what exactly is required for a workspace to qualify as ‘global’^{25,56}. Is it the number (and type) of consuming systems to which the workspace can broadcast that matters, or is it the kind of broadcasting that occurs within the workspace? Or are both of these considerations relevant to what counts as a ‘global workspace’? These questions need to be answered if we are to know what predictions GWTs make with respect to consciousness in, for example, infants, individuals with brain damage, split-brain patients, non-human animals and artificial intelligence systems.

Integrated information theory

IIT starts from a very different place than HOTs or GWTs by advancing a mathematical approach to characterising phenomenology. The theory starts by proposing axioms about the phenomenological character of conscious experiences (that is, properties that are taken to be self-evidently true and to apply to all possible forms of consciousness), and from these axioms deriving claims about the properties that any physical substrate of consciousness must satisfy. IIT then proposes that physical systems that instantiate these properties necessarily also instantiate consciousness (⁵⁷⁻⁶⁰, see **Fig. 3**).

Specifically, IIT proposes that consciousness should be understood in terms of ‘cause–effect power’ associated with irreducible maxima of integrated information generated by a physical system. Integrated information in turn is associated with the information theoretic quantity Φ (Phi), which measures — broadly speaking — how much information is generated by a system as a whole, compared with its parts considered independently. In IIT, consciousness is an intrinsic, fundamental property of a system, and is determined both by the nature of the causal mechanisms that compose it and by their state⁶⁰.

In contrast to HOTs and GWTs, IIT links consciousness primarily with posterior cortical areas (the so-called posterior ‘hot zone’ encompassing parietal, temporal and occipital areas) on the grounds that these areas exhibit neuroanatomical properties supposedly well suited for generating high levels of integrated information⁵⁹. Also in contrast to GWTs and HOTs, which associate consciousness with aspects of cortical information processing (that is, functional descriptions of what a system does), IIT does not refer to ‘information processing’ per se. Instead, it links consciousness to properties of the intrinsic cause–effect structure of a system: namely, to the causal power of a system to influence itself. According to IIT, any system that generates a non-zero maximum of (irreducible) integrated information is conscious, at least to some degree. Because of this, IIT would appear to imply that there already exist non-biological systems that are conscious⁶¹.

IIT is reasonably comprehensive, offering accounts of both global states and local states of consciousness (⁵⁹, see **Fig. 3**). Global states are associated with the quantity of irreducible integrated information generated by a system, as measured by Φ . IIT therefore encourages a unidimensional conception of global states, for it equates an

organism's level of consciousness with its value of Φ . The experiential character of local states can be understood in terms of 'conceptual structures', which IIT treats as 'shapes' in a high-dimensional space that is specified by the mechanistic cause-effect structure of the system. These shapes underpin (or are identical to) specific kinds of phenomenal character. For example, the spatial nature of visual experience has been related to the cause-effect structure specified by grid-like mechanisms present in early visual cortex⁶². The global unity of consciousness is explained by the integrated aspect of integrated information — its association with information generated by the 'whole' over and above that generated by the 'parts'. Finally, according to IIT, contents are conscious (rather than unconscious) when, and only when, they are incorporated into a cause-effect 'complex' (where a complex is a subset of the physical system that underpins a maximum of irreducible integrated information).

Returning to the binocular rivalry example, IIT explains why the subject reports experiencing a house (rather than a face) by appealing to the hypothesis that the complex underlying their report is associated with the conceptual structure corresponding to the content 'house' (rather than the content 'face'), and it explains the experiential contrast between seeing a house and seeing a face in terms of the 'shape' of the corresponding conceptual structure.

Although IIT provides a more comprehensive treatment of the various aspects of consciousness than most ToCs, it says comparatively little about how consciousness is related to other aspects of the mind, such as attention, learning and memory, nor has it yet focused on the relevance of embodiment and environmental embeddedness for consciousness (the latter also being a challenge for HOTs and GWTs). That said, IIT theorists have made initial steps towards addressing some of these challenges by, for example, developing measures of 'matching complexity' that track the shared information between an agent and its environment, and by formulating agent-based models in which agents that are able to engage effectively with their surroundings are found to exhibit increased amounts of integrated information⁶²⁻⁶⁴.

Reentry and predictive processing

Finally, we consider two general approaches to understanding consciousness that emphasize the importance of top-down signalling in shaping and enabling conscious perception. The first — reentry theories — are theories of consciousness as such, and associate conscious perception with top-down (recurrent, reentrant) signalling^{65,66}. The second group — PP theories — are not first-and-foremost ToCs but are more general accounts of brain (and body) function that can be used to formulate explanations and predictions regarding properties of consciousness⁶⁷.

Reentry theories are motivated by neurophysiological evidence revealing the importance of top-down signalling for conscious (usually visual) perception (for examples, see REFS⁶⁸⁻⁷⁰). In one prominent reentry theory — local recurrency theory — Lamme argues that localised recurrent or reentrant processing within perceptual cortices is sufficient to give rise to consciousness (given the integrity of other enabling factors, such as brainstem arousal), but that parietal and frontal regions might be

required for reporting the contents of perceptual experience or drawing on them for reasoning and decision-making^{65,71} (**Fig. 4**).

Broadly speaking, PP theories have two motivations. One motivation traces to considering the problem of perception as one of inference about the causes of sensory signals^{72,73}. The other — exemplified by the free energy principle⁷⁴ — appeals to fundamental constraints regarding control and regulation that apply to all systems that maintain their organisation over time (⁷⁵⁻⁷⁷, but see⁷⁸). Both lead to the notion that the brain implements a process of ‘prediction error minimisation’⁷⁹ that approximates Bayesian inference through the reciprocal exchange of (usually top-down) perceptual predictions and (usually bottom-up) prediction errors (⁸⁰, although see REF.⁸¹). Some expressions of PP, such as active inference, add the notion that sensory prediction errors can be minimized not only by updating predictions, but also by performing actions to bring about expected sensory data — thereby enabling a form of predictive control^{82,83}.

Although PP theories did not originate as ToCs, it has been suggested that they can furnish systematic correlations between neural mechanisms and phenomenological properties⁶⁷, where ‘systematic’ means having explanatory power guided by theoretical considerations, in contrast to mere empirical correlations as in the vanilla NCC approach. From this perspective, PP theories fulfil many of the desiderata for ToCs we outlined earlier, but may be best thought of as theories *for* consciousness science, rather than theories *of* consciousness, for there are many perspectives on precisely how PP relates to consciousness^{84,85}.

PP theories typically address local conscious states in terms of the content of top-down perceptual predictions^{73,79,86,87}: informally, perceptual content is given by the brain’s ‘best guess’ of the causes of its sensorium. The experiential character of a local state is specified by the nature of the perceptual predictions at play. For example, the phenomenology of ‘objecthood’ in vision may be accounted for by conditional predictions about the sensory consequences of actions^{87,88}, whereas the phenomenology of emotional states may be explained by the role of interoceptive predictions in regulating the organism’s physiological condition^{89,90}. The example of emotion highlights that PP theories, more than the others discussed here, encompass issues related to conscious selfhood^{13,77,91}.

PP can explain the distinction between conscious and unconscious states in terms of whether a mental state is part of a current ‘best guess’ (or optimal posterior) during perceptual inference. In the example of binocular rivalry, PP envisages two competing perceptual hypotheses (best-guesses), one of which ‘wins’, leading to perceptual dominance. Sensory signals from the alternative hypothesis accumulate as prediction error, which eventually lead to a perceptual transition, at which point sensory signals explained by the previously dominant best-guess now become a source of unexplained prediction error, and so the cycle repeats^{92,93}. (The experiential contrast between the house and the face would, as mentioned, be explained by properties of the corresponding perceptual predictions.) In those varieties of PP that emphasise active inference, a change in conscious content can only happen if perceptual belief updating

comes about through action (where action can be overt, such as a saccadic eye movement, or covert, such as a shift of attentional focus)^{76,94}.

PP theories do not generally deal with global states of consciousness, but it would be natural for them to appeal to the integrity of the relevant predictive processes in explaining distinctions among global states⁹⁵, in much the same way in which HOT accounts can appeal to the integrity of the relevant meta-representational machinery.

With respect to the functional dimensions of consciousness, both reentry and PP approaches provide clear treatments of the relationship between consciousness and attention. In local recurrency theory, as in GWTs, attention provides a selective boost to sensory signals so that they reach prefrontal and parietal regions, engaging conscious access⁷¹. In PP, attention is associated with the process of ‘precision weighting’, in which the estimated precision of sensory signals is modulated in ways intuitively equivalent to altering the signal-to-noise ratio or ‘gain’ of these signals^{74,96}; and in active inference, as mentioned, attentional sampling may be necessary for changes in conscious content^{76,94}.

Evaluating theories of consciousness

The range of data that have been appealed to in connection with the debate between rival theories of consciousness is vast, and we cannot hope to provide a full inventory of it here. Instead, we offer a selective overview of some current debates, highlighting the diversity of data that can be brought to bear on the evaluation of ToCs. (Some other empirical data generally used in support of each ToC are described in the legends of Figs 1–4.)

As a background point, it is important to recognize the holistic nature of theory-evaluation. Theories are not confirmed by a single finding; nor are they generally defeated by a single experiment. Instead, theory confirmation is typically an incremental process, in which one theory wins out over its rivals by providing an account of the target phenomenon that explains a wide range of data and can be integrated with successful theories in neighbouring domains⁹⁷⁻⁹⁹.

One obvious source of constraints on a ToC is the structure of consciousness. Although a number of structural features have been discussed in connection with ToCs, one structural feature of particular utility for contrasting ToCs is the unity of consciousness — the fact that the experiences that a single agent has at a time seem always to occur as the components of a single complex experience, one that fully captures what it’s like to be that agent²⁵. Different ToCs take very different attitudes to the unity of consciousness. IIT places considerable emphasis on the unity of consciousness. It not only assumes that consciousness is always unified but also appeals to the claim that consciousness is necessarily unified to motivate the association of consciousness with (maxima of) irreducible integrated information. Although GWTs do not emphasize the unity of consciousness in the way that IIT does, the association of consciousness with broadcast within a functionally integrated workspace suggests that they too may have the resources to provide a plausible account of the unity of consciousness. Other ToCs, such as HOTs and reentry/PP theories, have a more

ambivalent relationship to the unity of consciousness, tending either to only gesture at an account of this property or to overlook it entirely. The contrast in attitudes among ToCs to the unity of consciousness is due, in part at least, to more fundamental disagreement over whether consciousness is (necessarily) unified. Although the unity of consciousness promises to provide an important constraint on ToCs, in order for this promise to be realized we need a better account of the respects in which consciousness is (necessarily) unified.

A second source of constraints is provided by neural data. For example, it is generally accepted that the cerebellum is neither necessary nor sufficient for consciousness. A ToC ought to account for this fact, and explain why the cerebellum is not implicated in consciousness. Some ToCs readily provide such an explanation — for example, IIT argues that the cerebellum is not implicated in consciousness because its architecture is poorly suited for generating high levels of integrated information⁵⁹. But explanations of this sort lend specific support to a theory only if the account provided is more plausible than the accounts that might be provided by its competitors, and it is currently an open question whether that condition is satisfied. For example, advocates of HOTs could argue that the cerebellum lacks the capacity to support meta-representations of the relevant kind; proponents of GWTs can make the case that the cerebellum does not implement a global workspace; and reentrant and PP theorists can point to the absence of rich recurrent signalling in the cerebellum⁶⁵.

Although it is generally accepted that a ToC should explain why the cerebellum is not implicated in consciousness, there are other kinds of neural data that are much more controversial from the point of view of ToCs. An important example is provided by the debate about the role of prefrontal ('front-of-brain') processes in consciousness.

Using a variety of experimental paradigms, many neuroimaging studies have found prefrontal engagement for conscious (versus unconscious) perception⁴⁸, based on both regional activity^{100,101} and on functional connectivity between frontal and other regions¹⁰². A small number of primate studies have also found that conscious contents can be decoded from prefrontal activity patterns during binocular rivalry, continuous flash suppression, and rapid serial presentation of visual stimuli¹⁰³⁻¹⁰⁵; see also REF.¹⁰⁶ for a more complex picture in which content-relevant information was decoded from a wide range of both activated and deactivated cortical regions during an object recognition task. Lesion evidence, and evidence from brain stimulation has also been used to argue that the prefrontal activity is crucially implicated in consciousness; see REF.³⁰ for a review.

Advocates of HOTs and GWTs appeal to these findings to support their accounts over competing theories. In response, advocates of IIT and reentry theories argue that the observed prefrontal activity is a (non-necessary) consequence of consciousness and is probably associated with cognitive access to the contents of consciousness and the ability to provide behavioural reports, rather than with conscious perception per se^{107,108} (but see REF.¹⁰⁹). Those who defend this 'back-of-brain' perspective argue that posterior cortical processes — encompassing parts of perceptual and parietal cortex and precuneus — suffice for perceptual experience, and that front-of-brain processes

are not necessary. This claim is supported by so-called ‘no report’ studies which have tended to find diminished prefrontal engagement when subjects do not provide explicit reports about their perceptions^{6,110} (but see REF.¹¹¹). ‘Back-of-brain’ advocates draw on positive evidence in favour of a tight coupling between posterior activity and consciousness. For example, one innovative study probing for conscious contents during sleep using a serial awakening paradigm found that activity in posterior cortical regions predicted whether an individual would report dream experience, across both REM and non-REM sleep stages¹¹² (but see REF.¹¹³). Finally, the ‘front-of-brain’ interpretation of decoding studies is open to challenge, for showing that conscious contents can be ‘read out’ from a particular area does not establish that the brain itself is ‘reading out’ those contents from that area in a way that constitutes a relevant kind of meta-representation or global broadcast.

Although some aspects of the front-of-brain versus back-of-brain debate do indeed concern neurobiological data —for example, opinions differ on where the anatomical boundaries of the prefrontal cortex lie^{107,109} — at its heart is a disagreement about the relationship between consciousness and cognitive access: is it reasonable to take the availability of content for verbal report and the direct control of behaviour as a proxy for consciousness, or should investigations into the brain basis of consciousness remain neutral as to how exactly consciousness and cognitive access are related¹¹⁴ (see also **Box 3**)? Debate about this question is reflected in the attitudes that different ToCs take to cognitive access. GWTs place cognitive access at the heart of their account of consciousness, suggesting not only that the contents of consciousness are always available for cognitive access, but also that the processes underlying cognitive access (namely, ignition and global broadcast) serve as the basis of conscious experience (see REF.¹¹¹ for a recent nuance on this view). Other theories, such as IIT and local recurrency accounts, deny a close relationship between consciousness and cognitive access, holding that mental states can be conscious without being available for the direct control of thought and action, and also that mental states could in principle be available for the direct control of thought and action without being conscious. Although higher-order approaches are not committed to any particular relationship between consciousness and cognitive access, in practice their advocates generally assume that the contents of consciousness will be cognitively accessible (for example, see REF.⁴⁶), although perhaps not vice-versa.

Perhaps the most powerful source of data for evaluating rival ToCs involves novel predictions. Many of the most significant events in the history of science have involved the confirmation of novel predictions¹¹⁵. For example, general relativity received strong support from the fact that it predicted both the advance of the perihelion of Mercury and the way in which starlight grazing the Sun’s surface would be deflected. If a ToC were to make confirmed novel predictions, then it would be strongly supported, especially when compared with theories that failed to make the relevant prediction, or made different and incompatible predictions.

Many of the novel predictions that contemporary ToCs make are difficult to test. For example, both the reentry and IIT accounts predict that posterior cortical activity

can support conscious experience without contribution from anterior areas, but at present we lack reliable methods to verify such claims, since verification relies on subjective report (or at least, executive control in some form), which in turn requires anterior cortical activity. More dramatically, IIT predicts that consciousness is widely distributed throughout nature, including in many non-biological systems, and might even occur in systems that are as simple as photodiodes and single atoms (although, interestingly, not in strictly feedforward neural networks, see⁶¹). This prediction runs counter to widely-held assumptions about the distribution of consciousness, but it cannot be sensibly evaluated in the absence of robust methods for detecting the presence of consciousness in such systems (**Box 3**).

In some cases, methodological advances may bring novel predictions within reach of testability. One striking prediction, arising from IIT, is that changes in neural structure could lead to changes in conscious experience even when these changes do not give rise to changes in neural activity¹¹⁶. For example, inactive neurons in visual cortex may contribute to visual experience, while inactivated neurons would not⁵⁹. This prediction arises because, in IIT, it is the cause–effect structure specified by neural mechanisms that matters for consciousness. This means that if one intervenes in neural mechanisms, so as to change the cause–effect structure, then consciousness can change even if the corresponding neural dynamics do not change — a prediction that is particularly counterintuitive in the case where dynamics are absent (that is, for inactive neurons). Hypotheses like this, which do not readily follow from the other theories discussed here, may be testable using precise interventional methods, such as optogenetics, in animal models of perceptual decision-making¹¹⁷.

A particularly fruitful avenue for evaluating rival ToCs focuses on the temporal profile of conscious (as opposed to unconscious) processing, as reflected for example by event-related potentials (ERPs) in electrophysiological recordings. Some theorists (for example, see REF.¹¹⁸) argue that conscious perception has an early (120–200ms) onset following stimulus presentation, appealing to evidence suggestive of a robust correlation between perceptual consciousness and early onset modality-specific negative-going ERPs — called awareness negativity responses — while questioning the reliability of previously discussed later-onset signatures, such as the P3b (a positive-going ERP observed at ~300ms after stimulus onset). The early negativity highlighted by Dembski and colleagues has been found in both vision and audition, leading them to argue that there is a generalized early-onset response that robustly indexes perceptual consciousness. Such data point in favour of IIT and local reentry accounts of consciousness (but see REF.¹¹⁹ for a later cross-modal signature of conscious perception). Other theorists^{120,121} argue in favour of a much later onset (roughly, 250–400 ms) for perceptual consciousness. Besides the debated P3b, late-onset accounts are motivated by various perceptual phenomena that appear to match this timescale, including the psychological refractory period, the attentional blink and postdictive effects — the latter being of particular interest in showing that a delayed cue can retrospectively trigger conscious perception¹²². Candidate late-onset neural signatures of conscious perception include long-distance information sharing and bifurcation

dynamics^{49,111}. Evidence in favour of late-onset accounts of perceptual consciousness generally supports higher-order and global workspace ToCs. The debate between ‘early-onset’ and ‘late-onset’ accounts of perceptual experience is likely to remain a central topic of discussion for the foreseeable future. Note that the issue of the temporal profile of conscious processing is distinct from both the perception of duration¹²³ and from the temporal characteristics of a conscious ‘moment’^{124,125}, both of which reflect aspects of conscious content that ought to be explained by a ToC.

Moving forward

At present, ToCs are generally used as ‘narrative structures’ within the science of consciousness. Although they inform the interpretation of neural and behavioural data, it is still rare for a study to be designed with questions of theory validation in mind¹²⁶. Although there is nothing wrong with employing theories in this manner, future progress will depend on experiments that enable ToCs to be tested and disambiguated. We conclude our review by identifying three issues that need to be addressed for a mature regimen of theory-testing to flourish in consciousness science.

First, ToCs need to be developed with precision, for theories that appeal only to vague and imprecise constructs can generate only vague and imprecise predictions. For example, HOTs and PP and reentry theories need to specify the kinds of meta-representations, reentrant or predictive processes that are distinctive of consciousness or of specific aspects of consciousness; IIT needs to make precise its implications for the functional profile of consciousness and the impact of the environment and embodiment on consciousness; and GWTs need to provide a principled account of which workspaces qualify as ‘global’ in the relevant sense.

A promising approach here is to use computational models to bring mechanistic specificity to ToCs that may have been formulated in relatively abstract or conceptual terms. In addition to grounding the generation of fine-grained predictions, such models might also provide a shared language in which the relative merits of rival ToCs can be compared, which can be especially useful for comparing ToCs originating from different starting points. For example, computational models could reveal shared principles of top-down signalling among HOTs and re-entry and PP theories, while clarifying the distinctions between meta-representation (for example, see REF.³⁵) and global broadcast (for example, see REF.^{127,128}) that separate HOTs from GWTs¹²⁹. The development of computational models might also allow contrasts between ToCs to be reframed in terms of (potentially distributed) processes rather than, as is currently popular, in terms of broad neuroanatomical regions (for example, as in the debate between ‘front-of-brain’ and ‘back-of-brain’ theorists¹¹¹). A key challenge for the computational approach is to develop models that do not merely account for the functional features of consciousness but also account for its phenomenological properties — a challenge that can be described by the general labels of ‘computational phenomenology’ and ‘computational neurophenomenology’ (for examples, see REFS^{37,130}). This brings the additional challenge of how to validate, or disambiguate between, computational models using phenomenological data (for example, see

REF.¹³¹); a challenge which can be met, at least in part, by collecting subjective reports at the appropriate levels of phenomenological granularity (**Box 3**).

In addition to being made more specific, ToCs also need to be made more comprehensive. For the most part, ToCs have tended to focus on particular kinds of local states (perceptual experiences, with an emphasis on vision), on particular kinds of global states (ordinary waking awareness), and on particular kinds of conscious creatures (adult human beings). Although there are good reasons why theorists have tended to focus on a restricted class of conscious states and creatures — experimental accessibility being an important factor — a fully comprehensive ToC must do justice to the rich diversity of consciousness. With respect to local states, ToCs must go beyond perception and account also, for example, for affect, temporality, volition and thought. With respect to global states, ToCs must go beyond ordinary wakefulness and account also for the distinctive modes of consciousness associated with, for example, dreaming, meditation, disorders of consciousness and the psychedelic state. With respect to conscious creatures, ToCs must go beyond adult experience and address questions regarding consciousness in human infants, non-human animals, and even artificial systems. Although there is nothing wrong with ToCs that have a restricted focus, theories that provide a more comprehensive account of consciousness have obvious advantages over those that do not, especially if they can identify explanatory connections between different aspects of consciousness.

The third issue to be addressed is the measurement problem: how can trustworthy measures of consciousness be identified¹³². Solving this problem is crucial, for detailed and comprehensive ToCs are unlikely to be of much use unless we also have the capacity to verify their predictions. It is useful to distinguish two (closely-related) versions of the measurement problem. The first concerns the detection of conscious contents. Here the primary challenge is to identify ways of distinguishing conscious from unconscious mental states that do not make controversial assumptions about the functional profile of consciousness (such as that conscious contents must be reportable or otherwise available for high-level cognitive control)^{114,133,134}. The other version of the measurement problem focuses not on contents but on creatures. The questions here include how we might determine the distribution of consciousness in the animal kingdom¹³⁵; whether certain classes of cerebral organoids¹³⁶ or artificial intelligence systems¹³⁵⁻¹³⁷ are conscious; when consciousness first emerges in ontogenesis¹³⁸; and when it is retained in the context of traumatic brain injury¹³⁹. Here too the challenge is to develop ways of measuring consciousness that avoid controversial assumptions about its functional profile (**Box 3**).

Of course, the above challenges are already being addressed, to varying extents, by consciousness researchers. These efforts are now complemented by initiatives such as the adversarial collaboration model, which is encouraging proponents of ToCs to devise experiments with the specific goal of differentiating between alternative ToCs¹⁷. Consciousness remains scientifically controversial, yet there is every reason to think that the iterative development, testing and comparison of ToCs will lead to a much deeper understanding of this most central of mysteries.

Acknowledgements

A.K.S. and T.B. are Fellows in the CIFAR Program on Brain, Mind, and Consciousness. A.K.S. is additionally grateful to the European Research Council (Advanced Investigator Grant 101019254) and the Dr. Mortimer and Theresa Sackler Foundation.

Author contributions

The authors both contributed to all aspects of preparing the article.

Competing interests

The authors declare no competing interests.

Glossary

Neural correlates of consciousness (NCC): The minimal set of neural events jointly sufficient for a conscious state.

Adversarial collaboration: A research project in which proponents of different theories together design an experiment to distinguish their preferred theories, and agree in advance about how the outcome will favour one theory over the other(s).

Global state (of consciousness): an organism's overall state of consciousness, usually linked to arousal and behavioural responsiveness, and associated with 'level' of consciousness.

Local state (of consciousness): a particular conscious mental state, such as a conscious perception, emotion, or thought. Local states are also often called conscious contents.

Phenomenal character (of a conscious state): the experiential nature of a local state, such as the 'redness' of red, or the pain of a toothache – sometimes also called qualia.

***Meta-representation:** a mental representation that has as its target another mental representation

Cognitive access: a functional property whereby a mental state has access to a wide range of cognitive processes, usually including verbal and/or behavioural report.

Φ (integrated information): the amount of information specified by a system that is irreducible to that specified by its parts. There are many variations of Φ , each calculated differently and making different assumptions.

Posterior hot zone: a range of brain regions towards the rear of the cortex, including parietal, temporal, and occipital areas, as well as regions such as the precuneus.

***Complex:** in IIT, a subset of a physical system that underpins a maximum of irreducible integrated information.

Binocular rivalry: a phenomenon in which different images are presented to each eye, and conscious perception alternates between the two images.

***Interoceptive predictions:** predictions about the causes of sensory signals originating from within the body (interoception refers to perception of the body 'from within')

***Unity of consciousness:** the fact that the experiences that a single agent has at a time seem always to occur as the components of a single complex experience

Computational (neuro)phenomenology: the use of computational models to account for the phenomenal character of a conscious state in terms of (neural) mechanisms.

***The measurement problem(s):** the problems of (a) identifying whether a particular mental state is conscious, or (b) determining whether an organism or other system has the capacity to be conscious.

***Cerebral (brain) organoid:** laboratory-grown neural structure that self-organises into a system with cellular and network features resembling aspects of the developing human brain.

No-report paradigm: behavioural experiment in which participants do not provide subjective (verbal, behavioural) reports.

Explanatory gap intuition: the intuition that there is no prospect of a fully satisfying explanation of consciousness in physical, mechanistic terms.

Table 1 | A selection of theories of consciousness that are either neurobiological or potentially expressible in neurobiological terms.

Theory	Primary claim	Key references
Higher-order thought theory	Consciousness depends on meta-representations of lower-order mental states	31,46
Self-organising meta-representational theory	Consciousness is the brain's (meta-representational) theory about itself	34,140
Attended intermediate representation theory	Consciousness depends on the attentional amplification of intermediate level representations	141,142
Neuronal global workspace theory	Consciousness depends on ignition and broadcast within a neuronal global workspace where frontoparietal cortical regions play a central, hub-like role	47-49
Integrated information theory	Consciousness is identical to the cause-effect structure of a physical substrate that specifies a maximum of irreducible integrated information.	57,59,60
Information closure theory	Consciousness depends on non-trivial information closure with respect to an environment at particular coarse-grained scales	143
Dynamic core theory	Consciousness depends on a functional cluster of neural activity combining high levels of dynamical integration and differentiation	144
Neural Darwinism	Consciousness depends on reentrant interactions reflecting a history of value-dependent learning events shaped by selectionist principles	145,146
Local recurrency	Consciousness depends on local recurrent or reentrant cortical processing and promotes learning	65,71
Predictive processing	Perception depends on predictive inference of the causes of sensory signals; provides a framework for systematically mapping neural mechanisms to aspects of consciousness	67,73,79
Neurorepresentationalism	Consciousness depends on multi-level neurally-encoded predictive representations	84
Active inference	While views vary, in one version consciousness depends on temporally and counterfactually deep inference about self-generated actions	76; see also ⁹¹

Beast machine theory	Consciousness is grounded in allostatic control-oriented predictive inference	13,75,77; see also ⁹⁰
Neural subjective frame	Consciousness depends on neural maps of bodily state providing a first-person perspective	24
Self comes to mind theory	Consciousness depends on interactions between homeostatic routines and multilevel interoceptive maps, with affect and feeling at the core	23,147
Attention schema theory	Consciousness depends on a neurally-encoded model of the control of attention	148
Multiple drafts model	Consciousness depends on multiple (potentially inconsistent) representations rather than a single, unified representation that is available to a central system	149
Sensorimotor theory	Consciousness depends on mastery of the laws governing sensorimotor contingencies	88
Unlimited associative learning	Consciousness depends on a form of learning which enables an organism to link motivational value with stimuli or actions that are novel, compound, and non-reflex inducing	150
Dendritic integration theory	Consciousness depends on integration of top-down and bottom-up signalling at a cellular level	151
Electromagnetic field theory	Consciousness is identical to physically integrated, and causally active, information encoded in the brain's global electromagnetic field	152
Orchestrated objective reduction	Consciousness depends on quantum computations within microtubules inside neurons	18

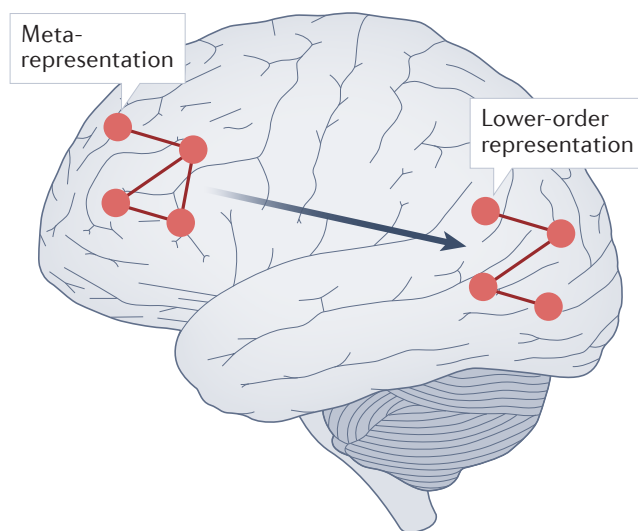


Fig. 1 | Higher-order theories. The core claim in higher-order theories (HOTs) of consciousness is that mental states are conscious in virtue of being the target of specific kinds of meta-representation. For example, lower-order representations of visual signals in posterior cortex would support conscious visual perception when targeted by the right kind of higher-order meta-representation. Supportive evidence for HOTs comes from studies implicating anterior cortical areas in conscious content, with an emphasis on prefrontal cortex — especially when performance is matched across conscious and non-conscious conditions^{30,100}. HOTs are also indirectly supported by lesion evidence linking metacognition to prefrontal areas¹⁵³. These theories are challenged by evidence suggesting that anterior areas are not involved in consciousness^{108,154}, perhaps instead being necessary only for enabling subjective report and executive control⁶. Figure adapted with permission from REF.⁴⁶.

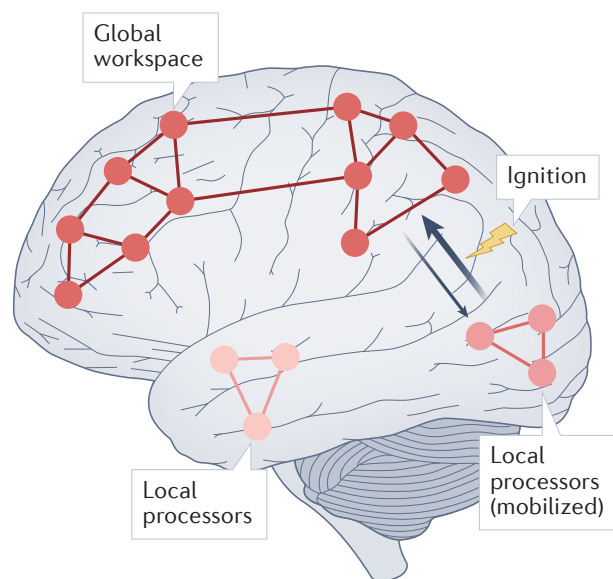


Fig. 2 | Global workspace theories. The core claim of global workspace theories (GWTs) of consciousness is that mental states are conscious when they are broadcast within a global workspace in which fronto-parietal networks play a central hub-like role. Activity in local processors (for example, sensory regions) becomes temporarily ‘mobilized’ into the workspace upon ignition¹⁵⁵. Empirical support for GWTs comes from studies associating consciousness with neuronal signatures of ignition and long-distance information sharing^{48,49,53,101}. Neural signatures of ignition are suggested by divergences of brain activity in anterior cortical regions at around 200ms to 300ms following stimulus onset, corresponding to trials with and without conscious perception^{48,101}, including in ‘no-report’ paradigms¹¹¹ (see also REF.¹⁵⁶). Such studies have been recently been extended to decoding: for example, activity patterns at around 300ms post-stimulus predicted subjective reports in ways that generalized across sensory modalities¹¹⁹. Signatures of long-distance information for conscious versus unconscious content have been identified using a range of methods^{49,102}. As with higher-order theories, GWTs are challenged by evidence that anterior regions might be involved in behavioural report rather than consciousness per se.

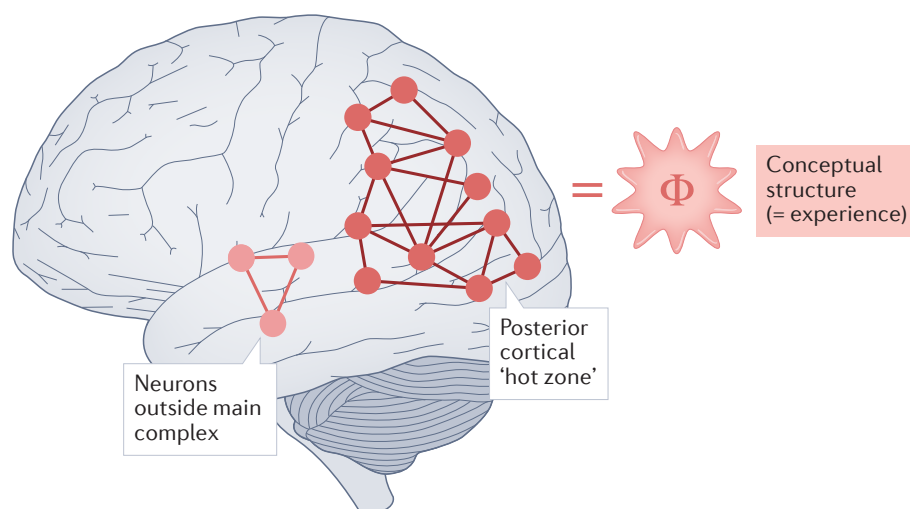


Fig. 3 | Integrated information theory. The core claim of integrated information theory (IIT) is that consciousness is identical to the cause-effect structure of a physical system that specifies a maximum of irreducible integrated information. The content of consciousness is associated with the form of the cause-effect structure, and the level of consciousness with its irreducibility, as measured by the quantity Φ ⁵⁹. Anatomically, IIT is associated with a posterior cortical ‘hot zone’. Empirical assessment of this core claim is challenging largely because Φ is infeasible to measure, except in simple model systems. Various proxies for Φ have been developed¹⁵⁷ and some show promise. Prominent among them is the ‘perturbation complexity index’ (PCI), which measures the algorithmic (Lempel-Ziv) complexity of the brain response to transcranial magnetic stimulation¹⁵⁸. Importantly, the PCI has diagnostic and prognostic value in tracking global states of consciousness in neurological patients¹⁵⁸. However, the PCI is not equivalent to Φ and correlations between the PCI and global states of consciousness are not incompatible with other theories of consciousness. Other evidence indirectly supportive of IIT comes from psychophysical studies suggesting that local changes in the strength of lateral connections within visual cortex can alter the structure of visual space¹¹⁶, and by evidence relating changes in global states to reduced functional diversity and integrative capacity in posterior cortical regions¹⁵⁹. IIT would be challenged by evidence which indicates that activity in anterior cortical regions is necessary for perceptual consciousness.

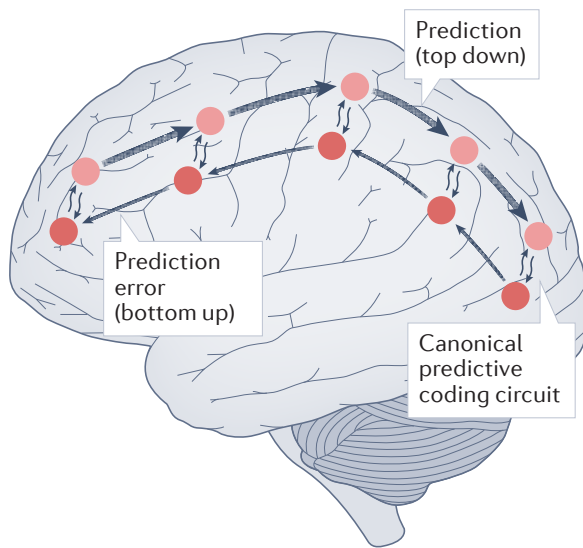


Fig. 4 | Reentry theory and predictive processing. The core claim of reentry theory and predictive processing (PP) is that conscious mental states are associated with top-down signalling (reentry, thick arrows) that, for PP, convey predictions about the causes of sensory signals (thin arrows signify bottom-up prediction errors), so that continuous minimisation of prediction errors implements an approximation to Bayesian inference. Conscious contents are specified – on most PP theories – by the content of the top-down predictions. Evidence in favour of these theories comes from studies linking top-down signalling with perceptual experience^{68-70,160}. In further support of PP, abundant evidence shows that expectations shape both the content of, and speed of access to, conscious perception¹⁶¹⁻¹⁶³ that some studies relate directly to top-down signalling¹⁶⁰. These theories would be challenged by evidence that top-down signalling or PP occurs in the absence of consciousness, or that changes in these processes do not affect conscious states.

Box 1 | Theories of consciousness and the ‘hard problem’

In the 1990s, David Chalmers famously distinguished between the ‘hard’ and ‘easy’ problems of consciousness¹⁶⁴. The easy problems are concerned with the functions and behaviours associated with consciousness, whereas the hard problem concerns the experiential (phenomenal, subjective) dimensions of consciousness. What makes the hard problem hard is the ‘explanatory gap’¹⁶⁵ — the intuition that there seems to be no prospect of a fully reductive explanation of experience in physical or functional terms.

Some theories of consciousness (for example, integrated information theory and certain versions of higher-order theory) address the hard problem directly. Other theories (for example, global workspace theories) focus on the functional and behavioural properties associated with consciousness; although they can be viewed as addressing the hard problem, this is not the primary goal of their proponents. A third strategy (adopted by some predictive processing theorists) aims to provide a framework in which various questions about the phenomenal properties of consciousness can be addressed, without attempting to account for the existence of phenomenology as such⁶⁷ – an approach sometimes called the ‘real problem’^{13,166}.

A critical question in this area is whether the hard problem is indeed a genuine challenge that ought to be addressed by a science of consciousness, or whether it ought to be dissolved rather than solved. Those who take the latter view often argue that the appearance of a distinctively hard problem derives from the peculiar features of the concepts (‘phenomenal concepts’) that we employ in representing our own conscious states^{167,168}. A related view is illusionism, according to which we do not actually have phenomenal states, but merely represent ourselves as having such states^{14,15}. Whatever the respective merits of these proposals, it seems likely that the grip of the hard problem may loosen as our capacity to explain, predict and control both phenomenological and functional properties of consciousness expands^{166,169}.

Box 2 | Other approaches: attention, learning and affect

The landscape of theories of consciousness includes a number of other theoretical approaches in addition to those surveyed in this review (**Table 1**). One approach focuses on attention. For example, Graziano's attention schema theory associates conscious perception with a model of the control of attention¹⁴⁸. Another attention-based theory of consciousness is the attended intermediate representational theory. First proposed in REF.¹⁴¹ and defended in detail in REF.¹⁴², this theory holds that consciousness occurs when intermediate-level perceptual representations gain access to attention.

Other theoretical approaches focus on learning. These include the proposal by Jablonka and Ginsburg that minimal consciousness is underpinned by a form of associative learning they term 'unlimited associative learning'. According to their proposal, this form of learning enables an organism to link motivational value with stimuli or actions that are novel, compound and non-reflex inducing¹⁵⁰. Other learning-based theories overlap with some theories we have already described, such as Cleeremans' version of higher-order theory^{34,140}, and Lamme's local recurrency account, which holds that recurrent signalling underpins consciousness in virtue of its role in learning⁶⁵. Learning-based theories are also closely related to 'selectionist' approaches, which ground consciousness in evolutionary-like dynamics within and between neuronal populations^{145,146}.

Affect-based theories emphasise the brain's role in physiological regulation as the basis for consciousness. These theories include Damasio's proposal that consciousness depends on hierarchically-nested representations of the organism's physiological condition^{147,170}, and proposals that mix an affect-based emphasis with predictive processing to ground conscious experiences in control-oriented interoceptive predictions^{13,77,90}. Some affect-based theories deny that cortical mechanisms are necessary for consciousness, instead locating the mechanisms of consciousness in the brainstem^{171,172} (although see REF.¹⁷³).

Box 3 | The measurement problem

To test a theory of consciousness, we need to be able to reliably detect both consciousness and its absence. At present, experimenters typically rely on a subject's introspective capacities, either directly or indirectly, to identify the states of consciousness experienced by that subject. However, this approach is problematic, for not only is the reliability of introspection questionable, there are many organisms or systems (for example, infants, individuals with brain damage and non-human animals) who might be conscious but are unable to produce introspective reports. Thus, there is a pressing need to identify non-introspective 'markers' or 'signatures' of consciousness.

A number of such indicators have been proposed in recent years. Some of these — such as the Perturbational Complexity Index (PCI)¹⁵⁸ — have been proposed as markers of consciousness as such, while others — such as the optokinetic nystagmus response¹⁷⁴ or distinctive bifurcations in neural dynamics¹¹¹ — have been proposed as markers of specific kinds of conscious contents. The former have been applied fruitfully to assessing global states of consciousness in individuals with brain injury¹⁷⁵ while the latter have been deployed in 'no report' studies of conscious content, in which overt behavioural reports are not made⁶. Whatever its focus, however, any proposed indicator of consciousness must be validated: we need to know that it is both sensitive and specific. While approaches to validation based on introspection have the problems mentioned above, theory-based approaches are also problematic. Because theories of consciousness are themselves contentious, it seems unlikely that appealing to theory-based considerations could provide the kind of intersubjective validation required for an objective marker of consciousness. Solving the measurement problem thus seems to require a method of validation that is neither based solely on introspection nor on theoretical considerations. The literature contains a number of proposals for addressing this problem^{114,176}, but none is uncontroversial^{177,178}.

- 1 Crick, F. & Koch, C. Towards a neurobiological theory of consciousness. *Seminars in the Neurosciences* **2**, 263-275 (1990).
- 2 Metzinger, T. (Ed.) *Neural correlates of consciousness: Empirical and conceptual questions*. MIT Press, Cambridge, MA (2000).
- 3 Koch, C., Massimini, M., Boly, M. & Tononi, G. Neural correlates of consciousness: progress and problems. *Nat Rev Neurosci* **17**, 307-321, doi:10.1038/nrn.2016.22 (2016).
- 4 de Graaf, T. A., Hsieh, P. J. & Sack, A. T. The 'correlates' in neural correlates of consciousness. *Neurosci Biobehav Rev* **36**, 191-197, doi:10.1016/j.neubiorev.2011.05.012 (2012).
- 5 Aru, J., Bachmann, T., Singer, W. & Melloni, L. Distilling the neural correlates of consciousness. *Neuroscience and biobehavioral reviews* **36**, 737-746, doi:10.1016/j.neubiorev.2011.12.003 (2012).
- 6 Tsuchiya, N., Wilke, M., Frassle, S. & Lamme, V. A. No-Report Paradigms: Extracting the True Neural Correlates of Consciousness. *Trends Cogn Sci* **19**, 757-770, doi:10.1016/j.tics.2015.10.002 (2015).
- 7 Klein, C., Hohwy, J. & Bayne, T. Explanation in the science of consciousness: From the neural correlates of consciousness (NCCs) to the difference-makers of consciousness (DMCs). *Philosophy and the Mind Sciences* **1** (2020).
- 8 Michel, M. *et al.* Opportunities and challenges for a maturing science of consciousness. *Nat Hum Behav* **3**, 104-107, doi:10.1038/s41562-019-0531-8 (2019).
- 9 Seth, A. K. Consciousness: The last 50 years (and the next). *Brain Neurosci Adv* **2**, 2398212818816019, doi:10.1177/2398212818816019 (2018).
- 10 Seth, A. K. Explanatory correlates of consciousness: Theoretical and computational challenges. *Cognitive Computation* **1**, 50-63 (2009).
- 11 Searle, J. *The rediscovery of the mind*. (MIT Press, 1992).
- 12 Varela, F. J. Neurophenomenology: A methodological remedy for the hard problem. *Journal of Consciousness Studies* **3**, 330-350 (1996).
- 13 Seth, A. K. *Being You: A New Science of Consciousness*. (Faber & Faber, 2021).
- 14 Dennett, D. C. Welcome to strong illusionism. *Journal of Consciousness Studies* **26**, 48-58 (2019).
- 15 Frankish, K. *Illusionism as a theory of consciousness*. (Imprint Academic, 2017).
- 16 Wiese, W. The science of consciousness does not need another theory, it needs a minimal unifying model. *Neurosci Conscious* **2020**, niaa013, doi:10.1093/nc/niaa013 (2020).
- 17 Melloni, L., Mudrik, L., Pitts, M. & Koch, C. Making the hard problem of consciousness easier. *Science* **372**, 911-912 (2021).
- Sets out how an adversarial collaboration is planning to arbitrate between integrated information and global workspace theories of consciousness.**
- 18 Hameroff, S. & Penrose, R. Consciousness in the universe: a review of the 'Orch OR' theory. *Phys Life Rev* **11**, 39-78, doi:10.1016/j.plrev.2013.08.002 (2014).
- 19 Chalmers, D. J. & McQueen, K. in *Quantum mechanics and consciousness* (ed S. Gao) (Oxford University Press, 2021).
- 20 Nagel, T. What is it like to be a bat? *Philosophical Review* **83**, 435-450 (1974).
- 21 Bayne, T., Hohwy, J. & Owen, A. M. Are There Levels of Consciousness? *Trends Cogn Sci* **20**, 405-413, doi:10.1016/j.tics.2016.03.009 (2016).

Challenges the common undimensional notion of ‘level of consciousness’, outlining an alternative, richer, multidimensional account.

- 22 Metzinger, T. *Being No-One*. (MIT Press, 2003).
- 23 Damasio, A. *Self comes to mind: Constructing the conscious brain*. (William Heinemann, 2010).
- 24 Park, H. D. & Tallon-Baudry, C. The neural subjective frame: from bodily signals to perceptual consciousness. *Philos Trans R Soc Lond B Biol Sci* **369**, 20130208, doi:10.1098/rstb.2013.0208 (2014).
- 25 Bayne, T. *The unity of consciousness*. (Oxford University Press, 2010).
- 26 Bayne, T. & Chalmers, D. J. in *The unity of consciousness: Binding, integration, and dissociation* (ed A. Cleeremans) 23-58 (Oxford University Press, 2003).
- 27 Cummins, R. Functional analysis. *The Journal of Philosophy* **72**, 741-765 (1975).
- 28 Blake, R., Brascamp, J. & Heeger, D. J. Can binocular rivalry reveal neural correlates of consciousness? *Philos Trans R Soc Lond B Biol Sci* **369**, 20130211, doi:10.1098/rstb.2013.0211 (2014).
- 29 Signorelli, C. M., Szczotka, J. & Prentner, R. Explanatory profiles of models of consciousness - towards a systematic classification. *Neurosci Conscious* **2021**, niab021, doi:10.1093/nc/niab021 (2021).
- 30 Lau, H. & Rosenthal, D. Empirical support for higher-order theories of conscious awareness. *Trends Cogn Sci* **15**, 365-373, doi:10.1016/j.tics.2011.05.009 S1364-6613(11)00105-7 [pii] (2011).
- Summary of empirical evidence favouring higher-order theories of consciousness.**
- 31 Rosenthal, D. *Consciousness and mind*. (Clarendon, 2005).
- 32 Brown, R. The HOROR theory of phenomenal consciousness. *Philos Stud* **172**, 1783-1794 (2015).
- 33 Cleeremans, A. Consciousness: the radical plasticity thesis. *Prog Brain Res* **168**, 19-33, doi:10.1016/S0079-6123(07)68003-0 (2008).
- 34 Cleeremans, A. *et al.* Learning to Be Conscious. *Trends Cogn Sci* **24**, 112-123, doi:10.1016/j.tics.2019.11.011 (2020).
- 35 Fleming, S. M. Awareness as inference in a higher-order state space. *Neurosci Conscious* **2020**, niz020, doi:10.1093/nc/niz020 (2020).
- 36 Lau, H. Consciousness, metacognition, and perceptual reality monitoring. *ArXiv*, doi:doi:10.31234/osf.io/ckbyf (2020).
- 37 Gershman, S. J. The generative adversarial brain. *Frontiers in artificial intelligence* **2**, doi:10.3389/frai.2019.00018 (2019).
- 38 Cohen, M. A., Dennett, D. C. & Kanwisher, N. What is the Bandwidth of Perceptual Experience? *Trends Cogn Sci* **20**, 324-335, doi:10.1016/j.tics.2016.03.006 (2016).
- 39 Haun, A. M., Tononi, G., Koch, C. & Tsuchiya, N. Are we underestimating the richness of visual experiences? *Neuroscience of Consciousness* **3**, 1-4 (2017).
- 40 Odegaard, B., Chang, M. Y., Lau, H. & Cheung, S. H. Inflation versus filling-in: why we feel we see more than we actually do in peripheral vision. *Philos Trans R Soc Lond B Biol Sci* **373**, doi:10.1098/rstb.2017.0345 (2018).
- 41 LeDoux, J. E. & Brown, R. A higher-order theory of emotional consciousness. *Proc Natl Acad Sci U S A* **114**, E2016-E2025, doi:10.1073/pnas.1619316114 (2017).
- 42 Morrison, J. Perceptual confidence. *Analytic Philosophy* **78**, 99-147 (2016).
- 43 Peters, M. A. K. Towards Characterizing the Canonical Computations Generating Phenomenal Experience. doi:doi:10.31234/osf.io/bqfr6 (2021).

- 44 Rosenthal, D. Consciousness and its function. *Neuropsychologia* **46** (2008).
- 45 Charles, L., Van Opstal, F., Marti, S. & Dehaene, S. Distinct brain mechanisms for
conscious versus subliminal error detection. *Neuroimage* **73**, 80-94,
doi:10.1016/j.neuroimage.2013.01.054 (2013).
- 46 Brown, R., Lau, H. & LeDoux, J. E. Understanding the Higher-Order Approach to
Consciousness. *Trends Cogn Sci* **23**, 754-768, doi:10.1016/j.tics.2019.06.009 (2019).
- 47 Baars, B. J. *A cognitive theory of consciousness*. (Cambridge University Press, 1988).
- 48 Dehaene, S. & Changeux, J. P. Experimental and theoretical approaches to conscious
processing. *Neuron* **70**, 200-227, doi:S0896-6273(11)00258-3 [pii]
10.1016/j.neuron.2011.03.018 (2011).
- 49 Mashour, G. A., Roelfsema, P., Changeux, J. P. & Dehaene, S. Conscious Processing
and the Global Neuronal Workspace Hypothesis. *Neuron* **105**, 776-798,
doi:10.1016/j.neuron.2020.01.026 (2020).
- Summary of the neuronal global workspace theory and its supporting evidence.**
- 50 Dehaene, S., Sergent, C. & Changeux, J. P. A neuronal network model linking
subjective reports and objective physiological data during conscious perception. *Proc
Natl Acad Sci U S A* **100**, 8520-8525, doi:10.1073/pnas.1332574100 (2003).
- 51 Naccache, L. Why and how access consciousness can account for phenomenal
consciousness. *Philos Trans R Soc Lond B Biol Sci* **373**, doi:10.1098/rstb.2017.0357
(2018).
- 52 Mashour, G. A. Cognitive unbinding: a neuroscientific paradigm of general
anesthesia and related states of unconsciousness. *Neurosci Biobehav Rev* **37**, 2751-
2759, doi:10.1016/j.neubiorev.2013.09.009 (2013).
- 53 Demertzi, A. *et al.* Human consciousness is supported by dynamic complex patterns
of brain signal coordination. *Sci Adv* **5**, eaat7603, doi:10.1126/sciadv.aat7603 (2019).
- A large empirical study of functional connectivity patterns across different global
states of consciousness, focusing on how these patterns relate to underlying
structural connectivity.**
- 54 Barttfeld, P. *et al.* Signature of consciousness in the dynamics of resting-state brain
activity. *Proc Natl Acad Sci U S A* **112**, 887-892, doi:10.1073/pnas.1418031112
(2015).
- 55 Uhrig, L. *et al.* Resting-state Dynamics as a Cortical Signature of Anesthesia in
Monkeys. *Anesthesiology* **129**, 942-958, doi:10.1097/ALN.0000000000002336
(2018).
- 56 Carruthers, P. *Human and animal minds: The consciousness questions laid to rest*.
(Oxford University Press, 2019).
- 57 Tononi, G. Consciousness as integrated information: a provisional manifesto. *Biol
Bull* **215**, 216-242 (2008).
- 58 Tononi, G. Integrated information theory of consciousness: an updated account.
Arch Ital Biol **150**, 293-329 (2012).
- 59 Tononi, G., Boly, M., Massimini, M. & Koch, C. Integrated information theory: from
consciousness to its physical substrate. *Nat Rev Neurosci* **17**, 450-461,
doi:10.1038/nrn.2016.44 (2016).
- An account the core claims and concepts of the integrated information theory of
consciousness.**

- 60 Oizumi, M., Albantakis, L. & Tononi, G. From the phenomenology to the mechanisms of consciousness: Integrated Information Theory 3.0. *PLoS Comput Biol* **10**, e1003588, doi:10.1371/journal.pcbi.1003588 (2014).
- 61 Tononi, G. & Koch, C. Consciousness: here, there and everywhere? *Philos Trans R Soc Lond B Biol Sci* **370**, doi:10.1098/rstb.2014.0167 (2015).
- 62 Haun, A. M. & Tononi, G. Why Does Space Feel the Way it Does? Towards a Principled Account of Spatial Experienc. *Entropy* **21**, 1160 (2019).
- 63 Albantakis, L., Hintze, A., Koch, C., Adami, C. & Tononi, G. Evolution of integrated causal structures in animats exposed to environments of increasing complexity. *PLoS Comput Biol* **10**, e1003966, doi:10.1371/journal.pcbi.1003966 (2014).
- 64 Marshall, W., Gomez-Ramirez, J. & Tononi, G. Integrated Information and State Differentiation. *Front Psychol* **7**, 926, doi:10.3389/fpsyg.2016.00926 (2016).
- 65 Lamme, V. A. Towards a true neural stance on consciousness. *Trends Cogn Sci* **10**, 494-501 (2006).
- 66 Lamme, V. A. & Roelfsema, P. R. The distinct modes of vision offered by feedforward and recurrent processing. *Trends Neurosci* **23**, 571-579 (2000).
- 67 Hohwy, J. & Seth, A. K. Predictive processing as a systematic basis for identifying the neural correlates of consciousness. *Philosophy and the Mind Sciences* **1**, 3 (2020).
- 68 Lamme, V. A., Super, H., Landman, R., Roelfsema, P. R. & Spekreijse, H. The role of primary visual cortex (V1) in visual awareness. *Vision Res* **40**, 1507-1521, doi:10.1016/s0042-6989(99)00243-6 (2000).
- 69 Pascual-Leone, A. & Walsh, V. Fast backprojections from the motion to the primary visual area necessary for visual awareness. *Science* **292**, 510-512, doi:10.1126/science.1057099 (2001).
- An early study using transcranial magnetic stimulation to reveal a role for reentrant activity in conscious visual perception in humans.**
- 70 Boehler, C. N., Schoenfeld, M. A., Heinze, H. J. & Hopf, J. M. Rapid recurrent processing gates awareness in primary visual cortex. *Proc Natl Acad Sci U S A* **105**, 8742-8747, doi:10.1073/pnas.0801999105 (2008).
- 71 Lamme, V. A. How neuroscience will change our view on consciousness. *Cogn Neurosci* **1**, 204-220, doi:10.1080/17588921003731586 (2010).
- 72 von Helmholtz, H. *Handuch der physiologik Optik*. (Voss, 1867).
- 73 Clark, A. Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behav Brain Sci* **36**, 181-204, doi:10.1017/S0140525X12000477 (2013).
- A classic exposition of predictive processing and its relevance for perception, cognition, and action.**
- 74 Friston, K. J. The free-energy principle: a unified brain theory? *Nat Rev Neurosci* **11**, 127-138, doi:nrn2787 [pii] 10.1038/nrn2787 (2010).
- 75 Seth, A. K. in *Open MIND* (eds J. M. Windt & T. Metzinger) 35(T) (MIND Group, 2015).
- 76 Friston, K. J. Am I Self-Conscious? (Or Does Self-Organization Entail Self-Consciousness?). *Front Psychol* **9**, 579, doi:10.3389/fpsyg.2018.00579 (2018).
- 77 Seth, A. K. & Tsakiris, M. Being a Beast Machine: The Somatic Basis of Selfhood. *Trends Cogn Sci* **22**, 969-981, doi:10.1016/j.tics.2018.08.008 (2018).
- 78 Bruineberg, J., Dolega, K., Dewhurst, J. & Baltieri, M. The Emperor's New Markov Blankets *Behavioral and Brain Sciences* ((in press)).

- 79 Hohwy, J. *The Predictive Mind*. (Oxford University Press, 2013).
- 80 Rao, R. P. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat Neurosci* **2**, 79-87, doi:10.1038/4580 (1999).
- 81 Teufel, C. & Fletcher, P. C. Forms of prediction in the nervous system. *Nat Rev Neurosci* **21**, 231-242, doi:10.1038/s41583-020-0275-5 (2020).
- 82 Friston, K. J., Daunizeau, J., Kilner, J. & Kiebel, S. J. Action and behavior: a free-energy formulation. *Biological Cybernetics* **102**, 227-260, doi:10.1007/s00422-010-0364-z (2010).
- 83 Parr, T. & Friston, K. J. Generalised free energy and active inference. *Biol Cybern* **113**, 495-513, doi:10.1007/s00422-019-00805-w (2019).
- 84 Pennartz, C. M. A. Consciousness, Representation, Action: The Importance of Being Goal-Directed. *Trends Cogn Sci* **22**, 137-153, doi:10.1016/j.tics.2017.10.006 (2018).
- 85 Williford, K., Bennequin, D., Friston, K. & Rudrauf, D. The Projective Consciousness Model and Phenomenal Selfhood. *Front Psychol* **9**, 2571, doi:10.3389/fpsyg.2018.02571 (2018).
- 86 Hohwy, J. New directions in predictive processing. *Mind and Language* (2020).
- 87 Seth, A. K. A predictive processing theory of sensorimotor contingencies: Explaining the puzzle of perceptual presence and its absence in synesthesia. *Cogn Neurosci* **5**, 97-118, doi:10.1080/17588928.2013.877880 (2014).
- 88 O'Regan, J. K. & Noë, A. A sensorimotor account of vision and visual consciousness. *Behav Brain Sci* **24**, 939-973; discussion 973-1031 (2001).
- The primary description of the sensorimotor theory of consciousness, which argues that conscious perception is intimately related to action.**
- 89 Seth, A. K. Interoceptive inference, emotion, and the embodied self. *Trends Cogn Sci* **17**, 565-573, doi:10.1016/j.tics.2013.09.007 (2013).
- A theoretical application of predictive processing to interoception and physiological regulation, relating this to experiences of emotion and selfhood.**
- 90 Barrett, L. F. The theory of constructed emotion: an active inference account of interoception and categorization. *Soc Cogn Affect Neurosci* **12**, 1833, doi:10.1093/scan/nsx060 (2017).
- 91 Solms, M. The Hard Problem of Consciousness and the Free Energy Principle. *Front Psychol* **9**, 2714, doi:10.3389/fpsyg.2018.02714 (2018).
- 92 Hohwy, J., Roepstorff, A. & Friston, K. Predictive coding explains binocular rivalry: an epistemological review. *Cognition* **108**, 687-701, doi:S0010-0277(08)00132-7 [pii] 10.1016/j.cognition.2008.05.010 (2008).
- 93 Parr, T., Corcoran, A. W., Friston, K. J. & Hohwy, J. Perceptual awareness and active inference. *Neurosci Conscious* **2019**, niz012, doi:10.1093/nc/niz012 (2019).
- 94 Friston, K. J., FitzGerald, T., Rigoli, F., Schwartenbeck, P. & Pezzulo, G. Active Inference: A Process Theory. *Neural Comput* **29**, 1-49, doi:10.1162/NECO_a_00912 (2017).
- 95 Boly, M. *et al.* Preserved feedforward but impaired top-down processes in the vegetative state. *Science* **332**, 858-862, doi:10.1126/science.1202043 (2011).
- A neuroimaging study using dynamic causal modelling to show that loss of consciousness in the vegetative state is associated with impaired top-down connectivity from frontal to temporal cortices.**

- 96 Parr, T. & Friston, K. J. Working memory, attention, and salience in active inference. *Scientific reports* **7**, 14678, doi:10.1038/s41598-017-15249-0 (2017).
- 97 Chalmers, A. *What is this thing called science?*, (Queensland University Press, 2013).
- 98 Godfrey-Smith, P. G. *Theory and reality: An introduction to the philosophy of science*. 2nd edn, (University of Chicago Press, 2021).
- 99 Lipton, P. *Inference to the best explanation*. (Routledge, 2004).
- 100 Lau, H. & Passingham, R. E. Relative blindsight in normal observers and the neural correlate of visual consciousness. *Proc Natl Acad Sci U S A* **103**, 18763-18768 (2006).
An empirical study comparing conscious and unconscious visual perception in humans, controlling for performance, and revealing differences in prefrontal activation.
- 101 van Vugt, B. *et al.* The threshold for conscious report: Signal loss and response bias in visual and frontal cortex. *Science* **360**, 537-542, doi:10.1126/science.aar7186 (2018).
An empirical study which tracked the time course of neural signals in primate frontal cortex, showing that perceived stimuli elicited sustained activity, when compared to non-perceived stimuli.
- 102 Gaillard, R. *et al.* Converging intracranial markers of conscious access. *PLoS Biol* **7**, e61, doi:10.1371/journal.pbio.1000061 (2009).
- 103 Panagiotaropoulos, T. I., Deco, G., Kapoor, V. & Logothetis, N. K. Neuronal discharges and gamma oscillations explicitly reflect visual consciousness in the lateral prefrontal cortex. *Neuron* **74**, 924-935, doi:10.1016/j.neuron.2012.04.013 (2012).
- 104 Kapoor, V. *et al.* Decoding the contents of consciousness from prefrontal ensembles. (2020).
- 105 Bellet, J. *et al.* Decoding rapidly presented visual stimuli from prefrontal ensembles without report nor post-perceptual processing. *Neurosci Conscious* **2022**, niac005, doi:10.1093/nc/niac005 (2022).
- 106 Levinson, M., Podvalny, E., Baete, S. H. & He, B. J. Cortical and subcortical signatures of conscious object recognition. *Nature communications* **12**, 2930, doi:10.1038/s41467-021-23266-x (2021).
- 107 Boly, M. *et al.* 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, doi:10.1523/JNEUROSCI.3218-16.2017 (2017).
- 108 Raccach, O., Block, N. & Fox, K. C. R. Does the Prefrontal Cortex Play an Essential Role in Consciousness? Insights from Intracranial Electrical Stimulation of the Human Brain. *J Neurosci* **41**, 2076-2087, doi:10.1523/JNEUROSCI.1141-20.2020 (2021).
- 109 Odegaard, B., Knight, R. T. & Lau, H. Should a Few Null Findings Falsify Prefrontal Theories of Conscious Perception? *J Neurosci* **37**, 9593-9602, doi:10.1523/JNEUROSCI.3217-16.2017 (2017).
- 110 Brascamp, J., Blake, R. & Knapen, T. Negligible fronto-parietal BOLD activity accompanying unreportable switches in bistable perception. *Nat Neurosci*, doi:10.1038/nn.4130 (2015).
An empirical 'no-report' study showing that front-parietal activity did not track switches in perceptual dominance when subjective reports were not required.
- 111 Sergent, C. *et al.* Bifurcation in brain dynamics reveals a signature of conscious processing independent of report. *Nature communications* **12**, 1149, doi:10.1038/s41467-021-21393-z (2021).

- 112 Siclari, F. *et al.* The neural correlates of dreaming. *Nat Neurosci* **20**, 872-878, doi:10.1038/nn.4545 (2017).
- 113 Wong, W. *et al.* The Dream Catcher experiment: blinded analyses failed to detect markers of dreaming consciousness in EEG spectral power. *Neurosci Conscious* **2020**, niaa006, doi:10.1093/nc/niaa006 (2020).
- 114 Block, N. Consciousness, accessibility, and the mesh between psychology and neuroscience. *Behavioral and Brain Sciences* **30**, 481-548 (2007).
Argues that research in psychology and neuroscience shows that there is a real and not merely conceptual distinction between phenomenal consciousness (i.e., experience) and cognitive access to phenomenal consciousness
- 115 Musgrave, A. in *Relativism and realism in science* (ed R. Nola) 229-252 (Kluwer, 1988).
- 116 Song, C., Haun, A. M. & Tononi, G. Plasticity in the Structure of Visual Space. *eNeuro* **4**, doi:10.1523/ENEURO.0080-17.2017 (2017).
- 117 Marshel, J. H. *et al.* Cortical layer-specific critical dynamics triggering perception. *Science* **365**, doi:10.1126/science.aaw5202 (2019).
- 118 Dembski, C., Koch, C. & Pitts, M. Perceptual awareness negativity: a physiological correlate of sensory consciousness. *Trends Cogn Sci* **25**, 660-670, doi:10.1016/j.tics.2021.05.009 (2021).
- 119 Sanchez, G., Hartmann, T., Fusca, M., Demarchi, G. & Weisz, N. Decoding across sensory modalities reveals common supramodal signatures of conscious perception. *Proc Natl Acad Sci U S A* **117**, 7437-7446, doi:10.1073/pnas.1912584117 (2020).
- 120 Sergent, C. The offline stream of conscious representations. *Philos Trans R Soc Lond B Biol Sci* **373**, doi:10.1098/rstb.2017.0349 (2018).
- 121 Michel, M. & Doerig, A. A new empirical challenge for local theories of consciousness. *Mind and Language* (2021).
- 122 Sergent, C. *et al.* Cueing attention after the stimulus is gone can retrospectively trigger conscious perception. *Curr Biol* **23**, 150-155, doi:10.1016/j.cub.2012.11.047 (2013).
An empirical study revealing that conscious perception of a stimulus can be influenced by events happening (hundreds of milliseconds) after the stimulus appeared ('retroperception')
- 123 Roseboom, W. *et al.* Activity in perceptual classification networks as a basis for human subjective time perception. *Nature communications* **10**, 267, doi:10.1038/s41467-018-08194-7 (2019).
- 124 Kent, L. & Wittmann, M. Special Issue: Consciousness science and its theories Time consciousness: the missing link in theories of consciousness. *Neurosci Conscious* **2021**, niab011, doi:10.1093/nc/niab011 (2021).
- 125 Husserl, E. *Ideas: A general introduction to pure phenomenology*. (Collier Books, 1963).
- 126 Yaron, I., Melloni, L., Pitts, M. & Mudrik, L. The ConTraSt database for analyzing and comparing empirical studies of consciousness theories. *Nature Human Behavior*. doi: 10.1038/s41562-021-01284-5 (2022).
An online resource of empirical studies of consciousness, organised with respect to different theories of consciousness.

- 127 Joglekar, M. R., Mejias, J. F., Yang, G. R. & Wang, X. J. Inter-areal Balanced Amplification Enhances Signal Propagation in a Large-Scale Circuit Model of the Primate Cortex. *Neuron* **98**, 222-234 e228, doi:10.1016/j.neuron.2018.02.031 (2018).
- 128 VanRullen, R. & Kanai, R. Deep learning and the Global Workspace Theory. *Trends Neurosci*, doi:10.1016/j.tins.2021.04.005 (2021).
- 129 Shea, N. & Frith, C. D. The Global Workspace Needs Metacognition. *Trends Cogn Sci* **23**, 560-571, doi:10.1016/j.tics.2019.04.007 (2019).
- 130 Suzuki, K., Roseboom, W., Schwartzman, D. J. & Seth, A. K. A Deep-Dream Virtual Reality Platform for Studying Altered Perceptual Phenomenology. *Scientific reports* **7**, 15982, doi:10.1038/s41598-017-16316-2 (2017).
- 131 Vilas, M. G., Auksztulewicz, R. & Melloni, L. Active inference as a computational framework for consciousness. *Review of Philosophy and Psychology* (2021).
- 132 Browning, H. & Veit, W. The measurement problem in consciousness. *Philosophical Topics* **48**, 85-108 (2020).
- 133 Seth, A. K., Dienes, Z., Cleeremans, A., Overgaard, M. & Pessoa, L. Measuring consciousness: relating behavioural and neurophysiological approaches. *Trends Cogn Sci* **12**, 314-321 (2008).
- 134 Michel, M. Calibration in consciousness science. *Erkenntnis*, 1-22 (2021).
- 135 Birch, J., Schnell, A. K. & Clayton, N. S. Dimensions of Animal Consciousness. *Trends Cogn Sci* **24**, 789-801, doi:10.1016/j.tics.2020.07.007 (2020).
- 136 Bayne, T., Seth, A. K. & Massimini, M. Are There Islands of Awareness? *Trends Neurosci* **43**, 6-16, doi:10.1016/j.tins.2019.11.003 (2020).
- An examination of the possibility of consciousness in isolated neural systems such as brain organoids, disconnected cortical hemispheres, and ex cranio brains.**
- 137 Dehaene, S., Lau, H. & Kouider, S. What is consciousness, and could machines have it? *Science* **358**, 486-492, doi:10.1126/science.aan8871 (2017).
- 138 Hu, H., Cusack, R. & Naci, L. Typical and disrupted brain circuitry for conscious awareness in full-term and pre-term infants. *BioRxiv*, doi:<https://doi.org/10.1101/2021.07.19.452937> (2021).
- 139 Haugg, A. *et al.* Do Patients Thought to Lack Consciousness Retain the Capacity for Internal as Well as External Awareness? *Frontiers in neurology* **9**, 492, doi:10.3389/fneur.2018.00492 (2018).
- 140 Cleeremans, A. The Radical Plasticity Thesis: How the Brain Learns to be Conscious. *Front Psychol* **2**, 86, doi:10.3389/fpsyg.2011.00086 (2011).
- 141 Jackendoff, R. *Consciousness and the computational mind*. (MIT Press, 1987).
- 142 Prinz, J. *The conscious brain: How attention engenders experience*. (Oxford University Press, 2012).
- 143 Chang, A. Y. C., Biehl, M., Yu, Y. & Kanai, R. Information Closure Theory of Consciousness. *Frontiers in psychology* **11**, 1504, doi:10.3389/fpsyg.2020.01504 (2020).
- 144 Tononi, G. & Edelman, G. M. Consciousness and complexity. *Science* **282**, 1846-1851 (1998).
- An early proposal of how measures of neural complexity might relate to phenomenological properties of (all) conscious experiences**
- 145 Edelman, G. M. *Neural Darwinism: The Theory of Neuronal Group Selection*. (Basic Books, Inc., 1987).
- 146 Edelman, G. M. *The remembered present*. (Basic Books, 1989).

- 147 Damasio, A. *The feeling of what happens: Body and emotion in the making of consciousness*. (Harvest Books, 2000).
- 148 Graziano, M. S. A. The attention schema theory: A foundation for engineering artificial consciousness. *Frontiers in Robotics and AI* **4**, 60, doi:10.3389/frobt.2017.00060 (2017).
- 149 Dennett, D. C. *Consciousness Explained*. (Little, Brown, and London, 1991).
- 150 Ginsburg, S. & Jablonka, E. *The Evolution of the Sensitive Soul: Learning and the Origins of Consciousness*. (MIT Press, 2019).
- 151 Aru, J., Suzuki, M. & Larkum, M. E. Cellular Mechanisms of Conscious Processing. *Trends Cogn Sci* **24**, 814-825, doi:10.1016/j.tics.2020.07.006 (2020).
- 152 McFadden, J. Integrating information in the brain's EM field: the cemi field theory of consciousness. *Neurosci Conscious* **2020**, niaa016, doi:10.1093/nc/niaa016 (2020).
- 153 Fleming, S. M., Ryu, J., Golfinos, J. G. & Blackmon, K. E. Domain-specific impairment in metacognitive accuracy following anterior prefrontal lesions. *Brain* **137**, 2811-2822, doi:10.1093/brain/awu221 (2014).
- 154 Fox, K. C. R. *et al.* Intrinsic network architecture predicts the effects elicited by intracranial electrical stimulation of the human brain. *Nat Hum Behav* **4**, 1039-1052, doi:10.1038/s41562-020-0910-1 (2020).
- 155 Dehaene, S. & Naccache, L. Towards a cognitive neuroscience of consciousness: basic evidence and a workspace framework. *Cognition* **79**, 1-37 (2001).
- 156 Sergent, C., Baillet, S. & Dehaene, S. Timing of the brain events underlying access to consciousness during the attentional blink. *Nat Neurosci* **8**, 1391-1400, doi:10.1038/nn1549 (2005).
- 157 Mediano, P. A. M., Seth, A. K. & Barrett, A. B. Measuring integrated information: comparison of candidate measures in theory and simulation. *Entropy* **21**, 17, doi:<https://doi.org/10.3390/e21010017> (2019).
- 158 Casali, A. G. *et al.* A theoretically based index of consciousness independent of sensory processing and behavior. *Science translational medicine* **5**, 198ra105, doi:10.1126/scitranslmed.3006294 (2013).
- An empirical study showing that a measure of the complexity of the cortical response to transcranial magnetic stimulation distinguished between a range of global conscious states, including disorders of consciousness.**
- 159 Luppi, A. I. *et al.* Consciousness-specific dynamic interactions of brain integration and functional diversity. *Nature communications* **10**, 4616, doi:10.1038/s41467-019-12658-9 (2019).
- 160 Hardstone, R. *et al.* Long-term priors influence visual perception through recruitment of long-range feedback. *Nature communications* **12**, 6288, doi:10.1038/s41467-021-26544-w (2021).
- 161 de Lange, F. P., Heilbron, M. & Kok, P. How do expectations shape perception? *Trends Cogn Sci*, doi:<https://doi.org/10.1016/j.tics.2018.06.002> (2018).
- 162 Melloni, L., Schwiedrzik, C. M., Muller, N., Rodriguez, E. & Singer, W. Expectations change the signatures and timing of electrophysiological correlates of perceptual awareness. *J Neurosci* **31**, 1386-1396, doi:10.1523/JNEUROSCI.4570-10.2011 [pii] 10.1523/JNEUROSCI.4570-10.2011 (2011).
- An empirical study using a perceptual hysteresis paradigm to show that expectations enhance and accelerate conscious perception.**

- 163 Pinto, Y., van Gaal, S., de Lange, F. P., Lamme, V. A. & Seth, A. K. Expectations accelerate entry of visual stimuli into awareness. *J Vis* **15**, 13, doi:10.1167/15.8.13 (2015).
- 164 Chalmers, D. J. Facing up to the problem of consciousness. *Journal of Consciousness Studies* **23**, 200-219 (1995).
- The classic statement of the philosophical distinction between the ‘hard’ and ‘easy’ problems of consciousness**
- 165 Levine, J. Materialism and qualia: The explanatory gap. *Pacific Philosophical Quarterly* **64**, 354-361 (1983).
- 166 Seth, A. K. The real problem. *Aeon* (2016).
- 167 Balog, K. in *The Oxford Handbook of Philosophy of Mind* (eds A. Beckermann, B. P. McLaughlin, & S. Walter) (Oxford University Press, 2009).
- 168 Perry, J. *Knowledge, Possibility, and Consciousness*. (MIT Press, 2001).
- 169 Varela, F. J., Thompson, E. & Rosch, E. *The embodied mind: Cognitive science and human experience*. (MIT Press, 1993).
- 170 Carvalho, G. B. & Damasio, A. Interoception and the origin of feelings: A new synthesis. *Bioessays* **43**, e2000261, doi:10.1002/bies.202000261 (2021).
- 171 Solms, M. *The Hidden Spring: A Journey to the Source of Consciousness*. (Profile Books, 2021).
- 172 Merker, B. Consciousness without a cerebral cortex: a challenge for neuroscience and medicine. *Behav Brain Sci* **30**, 63-81; discussion 81-134, doi:10.1017/S0140525X07000891 (2007).
- 173 Parvizi, J. & Damasio, A. Consciousness and the brainstem. *Cognition* **79**, 135-160, doi:10.1016/s0010-0277(00)00127-x (2001).
- 174 Naber, M., Frassle, S. & Einhauser, W. Perceptual rivalry: reflexes reveal the gradual nature of visual awareness. *PLoS One* **6**, e20910, doi:10.1371/journal.pone.0020910 (2011).
- 175 Casarotto, S. *et al.* Stratification of unresponsive patients by an independently validated index of brain complexity. *Ann Neurol* **80**, 718-729, doi:10.1002/ana.24779 (2016).
- 176 Shea, N. & Bayne, T. The Vegetative State and the Science of Consciousness. *Br J Philos Sci* **61**, 459-484, doi:10.1093/bjps/axp046 (2010).
- 177 Birch, J. The search for invertebrate consciousness. *Noûs* (2020).
- 178 Phillips, I. The methodological puzzle of phenomenal consciousness. *Philos Trans R Soc Lond B Biol Sci* **373**, doi:10.1098/rstb.2017.0347 (2018).