

13 that robust perception and cognition can be modelled independently of any ontolog-
14 ical assumptions about the world in which an agent is embedded. Any agent-world
15 interaction can, in particular, also be represented as an agent-agent interaction.

16 **Keywords:** Active inference; Complex networks; Computation; Learning; Memory; Plan-
17 ning; Predictive coding; Self representation; Reference frame; Turing completeness

18 1 Introduction

19 It is a natural and near-universal assumption that the world objectively has the properties
20 and causal structure that we perceive it to have; to paraphrase Einstein’s famous remark
21 (*cf.* Mermin, 1985), we naturally assume that the moon is there whether anyone looks at it
22 or not. Both theoretical and empirical considerations, however, increasingly indicate that
23 this assumption is not correct. Beginning with the now-classic work of Aspect, Dalibard
24 and Roger (1982), numerous experiments by physicists have shown that neither photon
25 polarization nor electron spin obey local causal constraints; within the past year, all rec-
26 ognized loopholes in previous experiments along these lines have been closed (Hensen et
27 al., 2015; Shalm et al., 2015; Giustina et al., 2015). The trajectories followed by either
28 light (Jacques et al., 2007) or Helium atoms (Manning, Khakimov, Dall and Truscott,
29 2015) through an experimental apparatus have been shown to depend on choices made
30 by random-number generators after the particle has fully completed its transit of the ap-
31 paratus. Optical experiments have been performed in which the causal order of events
32 within the experimental apparatus is demonstrably indeterminate (Rubino et al., 2016).
33 As both the positions and momenta of large organic molecules have now been shown to ex-
34 hibit quantum superposition (Eibenberger et al., 2013), there is no longer any justification
35 for believing that the seemingly counter-intuitive behavior observed in these experiments
36 characterizes only atomic-scale phenomena. These and other results have increasingly led
37 physicists to conclude that the classical notion of an observer-independent “objective” real-
38 ity comprising spatially-bounded, time-persistent “ordinary objects” and well-defined local
39 causal processes must simply be abandoned (e.g. Jennings and Leifer, 2015; Wiseman,
40 2015).

41 These results in physics are complemented within perceptual psychology by computational
42 experiments using evolutionary game theory, which consistently show that organisms that
43 perceive and act in accord with the true causal structure of their environments will be
44 out-competed by organisms that perceive and act only in accord with arbitrarily-imposed,
45 organism-specific fitness functions (Mark, Marion and Hoffman, 2010; reviewed by Hoff-
46 man, Singh and Prakash, 2015). These results, together with theorems showing that an
47 organism’s perceptions and actions can display symmetries that the structure of the en-
48 vironment does not respect (Hoffman, Singh and Prakash, 2015; Prakash and Hoffman,
49 in review) and that organisms responsive only to fitness will out-complete organisms that
50 perceive the true structure of the environment in all but a measure-zero subset of environ-
51 ments (Prakash, Stephens, Hoffman, Singh and Fields, in review), motivate the interface

52 theory of perception (ITP), the claim that perceptual systems, in general, provide only an
53 organism-specific “user interface” to the world, not a veridical representation of its struc-
54 ture (Hoffman, Singh and Prakash, 2015; Hoffman, 2016). According to ITP, the perceived
55 world, with its space-time structure, objects and causal relations, is a virtual machine im-
56 plemented by the coupled dynamics of an organism and its environment. Like any other
57 virtual machine, the perceived world is merely an interpretative or semantic construct; its
58 structure and dynamics bear no law-like relation to the structure and dynamics of its im-
59 plementation (e.g. Cummins, 1977). In software systems, the absence of any requirement
60 for a law-like relation between the structure and dynamics of a virtual machine and the
61 structure and dynamics of its implementation allows hardware and often operating system
62 independence; essentially all contemporary software systems are implemented by hierar-
63 chies of virtual machines for this reason (e.g. Goldberg, 1974; Tanenbaum, 1976; Smith
64 and Nair, 2005). The ontological neutrality with which ITP regards the true structure of the
65 environment is, therefore, analogous to the ontological neutrality of a software application
66 that can run on any underlying hardware.

67 The evolutionary game simulations and theorems supporting ITP directly challenge the
68 widely-held belief that perception, and particularly human perception is *veridical*, i.e. that
69 it reveals the observer-independent objects, properties and causal structure of the world.
70 While this belief has been challenged before in the literature (e.g. by Koenderink, 2015), it
71 remains the dominate view by far among perceptual scientists. Marr (1982), for example,
72 held that humans “very definitely do compute explicit properties of the real visible surfaces
73 out there, and one interesting aspect of the evolution of visual systems is the gradual move-
74 ment toward the difficult task of representing progressively more objective aspects of the
75 visual world” (p. 340). Palmer (1999) similarly states, “vision is useful precisely because it
76 is so accurate ... we have what is called veridical perception ... perception that is consistent
77 with the actual state of affairs in the environment” (p. 6). Geisler and Diehl (2003) claim
78 that “much of human perception is veridical under natural conditions” (p. 397). Trivers
79 (2011) agrees that “our sensory systems are organized to give us a detailed and accurate
80 view of reality, exactly as we would expect if truth about the outside world helps us to
81 navigate it more effectively” (p. xxvi). Pizlo, Sawada and Steinman (2014) emphasize
82 that “veridicality is an essential characteristic of perception and cognition. It is absolutely
83 essential. *Perception and cognition without veridicality would be like physics without the*
84 *conservation laws.*” (p. 227; emphasis in original). The claim of ITP is, in contrast, that
85 objects, properties and causal structure as normally conceived are *observer-dependent rep-*
86 *resentations* that, like virtual-machine states in general, may bear no straightforward or
87 law-like relation to the actual structure or dynamics of the world. Evidence that specific
88 aspects of human perception are non-veridical, e.g. the narrowing and flattening of the
89 visual field observed by Koenderink, van Doorn and Todd (2009), the distortions of per-
90 spective observed by Pont et al. (2012), or the inferences of three-dimensional shapes from
91 motion patterns projectively inconsistent with such shapes observed by He, Feldman and
92 Singh (2015) provide *prima facie* evidence for ITP.

93 The implication of either ITP or quantum theory that the objects, properties and causal

94 relations that organisms perceive do not objectively exist as such raises an obvious challenge
95 for models of perception as an information-transfer process: the naïve-realist assumption
96 that perceptions of an object, property or causal process X are, in ordinary circumstances,
97 results of causal interactions with X cannot be sustained. Hoffman and Prakash (2014)
98 proposed to meet this challenge by developing a minimal, implementation-independent formal
99 framework for modelling perception and action analogous to Turing’s (1936) formal
100 model of computation. This “conscious agent” (CA) framework posits entities or systems
101 aware of their environments and acting in accordance with that awareness as its funda-
102 mental ontological assumption. The CA framework is a minimal refinement of previous
103 formal models of perception and perception-action cycles (Bennett, Hoffman and Prakash,
104 1989). Following Turing’s lead, the CA framework is intended not as a scientific or even
105 philosophical *theory* of conscious awareness, but rather as a minimal, universally-applicable
106 formal *model* of conscious perception and action. The universality claim made by Hoffman
107 and Prakash (2014) is analogous to the Church-Turing thesis of universality for the Turing
108 machine. Hoffman and Prakash (2014) showed that CAs may be combined to form larger,
109 more complex CAs and that the CA framework is Turing-equivalent and therefore univer-
110 sal as a representation of computation; this result is significantly elaborated upon in what
111 follows.

112 The present paper extends the work of Hoffman and Prakash (2014) by showing that the
113 CA framework provides a robust and intuitive representation of perceptual and cognitive
114 processes in the context of ITP. Anticipation, expectations and generative models of the
115 environment, in particular, emerge naturally in all but the simplest CA networks, providing
116 support for the claimed universality of the CA framework as a model of agent - world
117 interactions. We first define CAs and distinguish the *extrinsic* (external or “3rd person”)
118 perspective of a theorist describing a CA or network of CAs from the *intrinsic* (internal
119 or “1st person”) perspective of a particular CA. Consistency between these perspectives
120 is required by ITP; a CA cannot, in particular, be described as differentially responding
121 to structure in its environment that ITP forbids it from detecting. Such consistency can
122 be achieved by the “conscious realism” assumption (Hoffman and Prakash, 2014) that
123 the world in which CAs are embedded is composed entirely of CAs. We show that the
124 CA framework allows the incorporation of Bayesian inference from “images” to “scene
125 interpretations” as described by Hoffman and Singh (2012) and show that a CA can be
126 regarded as incorporating a “Markov blanket” as employed by Friston (2013) when this
127 is done. We analyze the behavior of the simplest networks of CAs in detail from the
128 extrinsic perspective, and discuss the formal structure and construction of larger, more
129 complex networks. We show that a concept of “fitness” for CAs emerges naturally within
130 the formalism, and that this concept corresponds to concepts of “centrality” already defined
131 within social-network theory. We then consider the fundamental question posed by ITP:
132 that of how non-veridical perception can be useful. We show that CAs can be constructed
133 that implement short- and long-term memory, categorization, active inference, goal-directed
134 attention, and case-based planning. Such complex CAs represent their world to themselves
135 as composed of “objects” that recur in their experience, and are capable of rational actions
136 with respect to such objects. This construction shows that specific ontological assumptions

137 about the world in which a cognitive agent is embedded, including the imposition of *a priori*
138 fitness functions, are unnecessary for the theoretical modelling of useful cognition. The non-
139 veridicality of perception implied by ITP need not, therefore, be regarded as negatively
140 impacting the behavior of an intelligent system in a complex, changing environment.

141 2 Conscious agents: Definition and interpretation

142 2.1 Definition of a CA

143 As noted, the CA framework is motivated by the hypothesis that agents of interest to
144 psychology are *aware* of the environments in which they act, even if this awareness is rudi-
145 mentary by typical human standards (Hoffman and Prakash, 2014). Our goal here is to
146 develop a minimal and fully-general formal model of perception, decision and action that
147 is applicable to any agent satisfying this hypothesis. Minimality and generality can be
148 achieved using a formalism based on measurable sets and Markovian kernels as described
149 below. This formalism allows us to explore the dynamics of multi-agent interactions (§3)
150 and the internal structures and dynamics, particularly of memory and attention systems,
151 that enable complex cognition (§4) constructively. We accordingly impose no *a priori* as-
152 sumptions regarding behavioral reportability or other criteria for inferring, from the outside,
153 that an agent is conscious *per se* or is aware of any particular stimulus; nor do we impose
154 any *a priori* distinction between conscious and unconscious states. Considering results such
155 as those reviewed by Boly, Sanders, Mashour and Laureys (2013), we indeed regard such
156 criteria and distinctions, at least as applied to living humans, as conceptually untrustwor-
157 thy and possibly incoherent. We thus treat awareness or consciousness as fundamental and
158 irreducible properties of agents, and ask, setting aside more philosophical concerns (but
159 see Hoffman and Prakash, 2014 for extensive discussion), what structural and dynamic
160 properties such agents can be expected to have.

161 We begin by defining the fundamental mathematical notions on which the CA framework
162 is based; we then interpret these notions in terms of perception, decision and action.

163 **Definition 1.** *Let $\langle B, \mathcal{B} \rangle$ and $\langle C, \mathcal{C} \rangle$ be measurable spaces. Equip the unit interval $[0, 1]$
164 with its Borel σ -algebra. We say that a function $K: B \times C \rightarrow [0, 1]$ is a **Markovian kernel**
165 **from B to C** if:*

166 (i) *For each measurable set $E \in \mathcal{C}$, the function $K(\cdot, E) : B \rightarrow [0, 1]$ enacted by $b \mapsto K(b, E)$
167 is a measurable function; and*

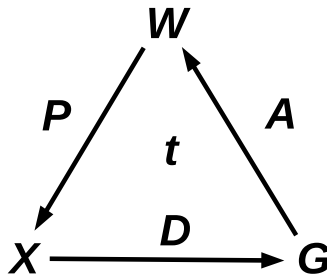
168 (ii) *For each $b \in B$, the function $K(b, \cdot)$ enacted by $F \mapsto K(b, F)$, $F \in \mathcal{C}$ is a probability
169 measure on C .*

170 In particular, if K is a Markovian kernel from B to C , then for any measurable $D \subset C$, the
171 function enacted by $x \mapsto K(x, D) \in [0, 1]$ assigns to each x in B a probability distribution

172 on C . When the spaces involved are finite, the Markov kernel can be represented as a
 173 matrix whose rows sum to unity.

174 We represent a CA as a labelled directed graph as shown in Fig. 1. This graph implies the
 175 development of a cyclic process, in which we can think of, e.g. the kernel $D : X \times G \rightarrow G$
 176 as follows: for each instantiation g_0 of G in the immediately previous cycle, and the current
 177 instantiation of $x \in X$, $D(x, g_0; \cdot)$ gives the probability distribution of the $g \in G$ instantiated
 178 at the next step. The other kernels A and P are interpreted similarly. Formally,

179 **Definition 2.** Let $\langle W, \mathcal{W} \rangle$, $\langle X, \mathcal{X} \rangle$ and $\langle G, \mathcal{G} \rangle$ be measurable spaces. Let P be a
 180 Markovian kernel $P : W \times X \rightarrow X$, D be a Markovian kernel $D : X \times G \rightarrow G$, and
 181 A be a Markovian kernel $A : G \times W \rightarrow W$. A **conscious agent (CA)** is a 7-tuple
 182 $[(X, \mathcal{X}), (G, \mathcal{G}), (W, \mathcal{W}), P, D, A, t]$, where t is a positive integer parameter.



183 *Fig. 1:* Representation of a CA as a labelled directed graph. W , X and G
 184 and measurable sets, P , D , and A are Markovian kernels, and t is an integer
 185 parameter.

186 Hoffman and Prakash (2014) defined a CA, given the measurable space $\langle W, \mathcal{W} \rangle$, as a
 187 6-tuple $[(X, \mathcal{X}), (G, \mathcal{G}), P, D, A, t]$ where $P : W \times X \rightarrow [0, 1]$, $D : X \times G \rightarrow [0, 1]$ and
 188 $A : G \times W \rightarrow [0, 1]$ are Markovian kernels and t is a positive integer parameter. Here
 189 we explicitly include $\langle W, \mathcal{W} \rangle$ in the definition of a CA. Following Hoffman, Singh and
 190 Prakash (2015) and Prakash and Hoffman (in review), we also explicitly allow the P , D ,
 191 and A kernels to depend on the elements of their respective target sets. Informally, for
 192 $x \in X$ and $g \in G$, for example, and any measurable $H \subset G$, the function enacted by
 193 $(x, g) \mapsto K(x, g, H)$ is real-valued and can be considered to be the regular conditional
 194 probability distribution $\text{Prob}(H|x, g)$ under appropriate conditions on the spaces involved
 195 (Parthasarathy, 2005). The difference in representational power between the more general,
 196 target-set dependent kernels specified here and the original, here termed “forgetful,” kernels
 197 of Hoffman and Prakash (2014) is discussed below.

198 We interpret elements of W as representing states of the “world,” making no particular
 199 ontological assumption about the elements or states of this world. We interpret elements of

200 X and G as representing possible conscious experiences and actions (strictly speaking, they
201 consist of formal *tokens* of possible conscious experiences and actions), respectively. The
202 kernels P, D and A represent perception, decision and action operators, where “perception”
203 includes *any* operation that changes the state of X , “decision” is any operation that changes
204 the state of G and “action” is any operation that changes the state of W . The set X is, in
205 particular, taken to represent all experiences regardless of modality; hence P incorporates
206 all perceptual modalities. The set G and kernel A are similarly regarded as multi-modal.
207 With this interpretation, perception can be viewed as an action performed by the world;
208 how these “actions” can be unpacked into the familiar bottom-up and top-down components
209 of perceptual experience is explored in detail in §4 below. The kernels P, D and A are taken
210 to act whenever the states of W, X or G , respectively, change. Both the decisions D and
211 the actions A of the CA are regarded as “freely chosen” in a way consistent with the
212 probabilities specified by D and A , as are the actions “by the world” represented by P ;
213 these operators are treated as stochastic in the general case to capture this freedom from
214 determination. The parameter t is a CA-specific proper time; t is regarded as “ticking”
215 and hence incrementing concurrently with the action of D , i.e. immediately following each
216 change in the state of X . No specific assumption is made about the contents of X ; in
217 particular, it is not assumed that X includes tokens representing the values of either t or
218 any elements of G . A CA need not, in other words, in general experience either time or its
219 own actions; explicitly enabling such experiences for a CA is discussed in §4.1 below.

220 It will be assumed in what follows that the contents of X and G can be considered to be
221 representations encoded by finite numbers of bits; for simplicity, all representations in X
222 or G will be assumed to be encoded, respectively, by the same numbers of bits. Hence X
223 and G can both be assigned a “resolution” with which they encode, respectively, inputs
224 from and outputs to W . It is, in this case, natural to regard D as operating in discrete
225 steps; for each previous instantiation of G , D maps one complete, fully-encoded element of
226 X to one complete, fully-encoded element of G . As the minimal size of a representation in
227 either X or G is one bit, the minimal action of D is a mapping of one bit to one bit. While
228 the CA framework as a whole is purely formal, we envision finite CAs to be amenable to
229 physical implementation. If any such physical implementation is assumed to be constrained
230 by currently accepted physics and the action of D is regarded as physically (as opposed
231 to logically) irreversible, the minimal energetic cost of executing D is given by Landauer’s
232 (1961; 1999) principle as $\ln 2 kT$, where k is Boltzmann’s constant and T is temperature in
233 degrees Kelvin. In this case, the minimal unit of t is given by $t = h/(\ln 2 kT)$, where h
234 is Planck’s constant. At $T \sim 310K$, physiological temperature, this value is $t \sim 100 fs$,
235 roughly the response time of rhodopsin and other photoreceptors (Wang et al., 1994). At
236 even the $50 ms$ timescale of visual short-term memory (Vogel, Woodman and Luck, 2006),
237 this minimal discrete time would appear continuous. As elaborated further below, however,
238 no general assumption about the coding capacities in bits of X or G are built into the CA
239 framework. What is to count, in a specific model, as an execution of D and hence an
240 incrementing of t is therefore left open, as it is in other general information-processing
241 paradigms such as the Turing machine.

242 Hoffman and Prakash (2014) explicitly proposed the “Conscious agent thesis: Every prop-
243 erty of consciousness can be represented by some property of a dynamical system of con-
244 scious agents” (p. 10), where the term “conscious agent” here refers to a CA as defined
245 above. As CAs are explicitly *formal models* of real conscious agents such as human be-
246 ings, the “properties of consciousness” with which this thesis is concerned are the *formal*
247 or computational properties of consciousness, e.g. the formal or computational properties
248 of recall or the control of attention, not their phenomenal properties. The conscious agent
249 thesis is intended as an empirical claim analogous to the Church-Turing thesis. Just as the
250 demonstration of a computational process not representable as a Turing machine computa-
251 tion would falsify the Church-Turing thesis, the demonstration of a conscious process, e.g.
252 a process of conscious recognition, inference or choice, not representable by the action of
253 a Markov kernel would falsify the conscious agent thesis. We offer in what follows both
254 theoretically-motivated reasons and empirical evidence to support the conscious agent the-
255 sis as an hypothesis. Whether the actual implementations of conscious processes in human
256 beings or other organisms can in fact be fully captured by a representation based on Markov
257 kernels remains an open question.

258 2.2 Extrinsic and intrinsic perspectives

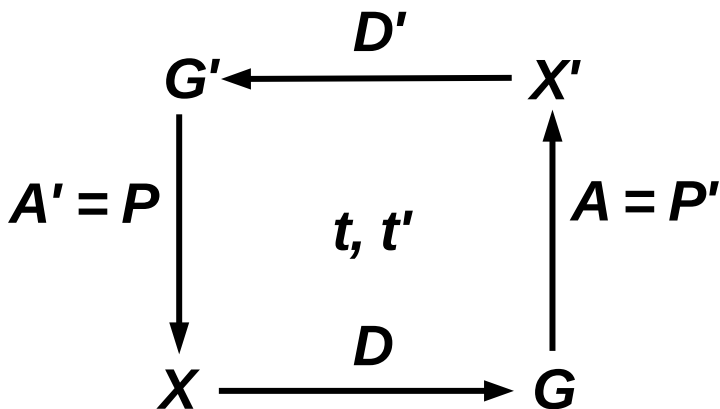
259 A central claim of ITP is that perceptual systems do not, in general, provide a veridical
260 representation of the structure of the world; in particular, “objects” and “causal relations”
261 appearing as experiences in X are in general not in any sense homomorphic to elements or
262 relationships between elements in W . This claim is, clearly, formulated from the extrinsic
263 perspective of a theorist able to examine the behavior of a CA “from the outside” and to
264 determine whether the kernel P is a homomorphism of W or not. The evolutionary game
265 theory experiments reported by Mark, Marion and Hoffman (2010) were conducted from
266 this perspective. As is widely but not always explicitly recognized, the extrinsic perspective
267 is of necessity an “as if” conceit; a theorist can at best construct a formal representation
268 of a CA and ask how the interaction represented by the $P - D - A$ cycle would unfold if it
269 had particular formal properties (e.g. Koenderink, 2014). The extrinsic perspective is, in
270 other words, a perspective of *stipulation*; it is not the perspective of any observer. For the
271 present purposes, the extrinsic perspective is simply the perspective from which the kernels
272 P , D and A may be formally specified.

273 The extrinsic perspective of the stipulating theorist contrasts with another relevant perspec-
274 tive, the intrinsic perspective of the CA itself. That every CA has an intrinsic perspective
275 is a consequence of the intended interpretation of CAs as *conscious* agents that experience
276 their worlds. Hence every CA is an observer, and the intrinsic perspective is the observer’s
277 perspective. The intrinsic perspective of a CA is most clearly formulated using the concept
278 of a “reduced CA” (RCA), a 4-tuple $[(X, \mathcal{X}), (G, \mathcal{G}), D, t]$. The RCA, together with a choice
279 of extrinsic elements W , A and P , is then what we have defined above as a CA. An RCA
280 can be viewed as both *embedded in* and *interacting with* the world represented by W . The
281 RCA freely chooses the action(s) to take - the element(s) of G to select - in response to

282 any experience $x \in X$; this choice is represented by the kernel D . The action A on W
 283 that the RCA is *capable* of taking is determined, in part, by the structure of W . Similarly,
 284 the action P with which W can affect the RCA is determined, in part, by the structure
 285 of the RCA. With this terminology, the central claim of ITP is that an RCA’s possible
 286 knowledge of W is completely specified by X ; the element(s) of X that are selected by P
 287 at any given t constitute the RCA’s entire experience of W at t . The structure and content
 288 of X completely specify, therefore, the intrinsic perspective of the RCA. In particular, ITP
 289 allows the RCA no independent access to the ontology of W ; consistency between intrinsic
 290 and extrinsic perspectives requires that no such access is attributed to any RCA from the
 291 latter perspective. An RCA does not, in particular, have access to the definitions of its
 292 own P , D or A kernels; hence an RCA has no way to determine whether any of them are
 293 homomorphisms. Similarly, an RCA has no access to the definitions of any other RCA’s P ,
 294 D or A kernels, or to any other RCA’s X or G . An RCA “knows” what currently appears
 295 as an experience in its own X but nothing else; as discussed in §4.1 below, for an RCA
 296 even to know what actions it has available or what actions it has taken in the past, these
 297 must be represented explicitly in X . Any structure attributed to W from the intrinsic
 298 perspective of an RCA is hypothetical in principle; such attributions of structure to W can
 299 be disconfirmed by continued observation, i.e. additional input to X , but can never be
 300 confirmed. In this sense, any RCA is in the epistemic position regarding W that Popper
 301 (1963) claims characterizes all of science.

302 From the intrinsic perspective, an immediate consequence of the ontological neutrality of
 303 ITP is that an RCA cannot determine, by observation, that the internal dynamics of its
 304 associated W is non-Markovian; hence it cannot distinguish W , as a source of experiences
 305 and a recipient of actions, from a second RCA. The RCA $[(X, \mathcal{X}), (G, \mathcal{G}), D, t]$, in partic-
 306 ular, cannot distinguish the interaction with W shown in Fig. 1 from an interaction with
 307 a second RCA $[(X', \mathcal{X}'), (G', \mathcal{G}'), D', t']$ as shown in Fig. 2. From the extrinsic perspective
 308 of a theorist, Fig. 2 can be obtained from Fig. 1 by interpreting the perception kernel P
 309 as representing actions by W on the RCA $[(X, \mathcal{X}), (G, \mathcal{G}), D, t]$ embedded within it. Each
 310 such action $P(w, \cdot)$ generates a probability distribution of experiences x in X . If an agent’s
 311 perceptions are to be regarded as actions on the agent by its world W , however, nothing
 312 prevents similarly regarding the agent’s actions on W as “perceptions” of W . If W both per-
 313 ceives and acts, it can itself be regarded as an agent, i.e. an RCA $[(X', \mathcal{X}'), (G', \mathcal{G}'), D', t']$,
 314 where the kernel D' represents W ’s internal dynamics. This symmetric interpretation of
 315 action and perception from the extrinsic perspective, with its concomitant interpretation
 316 of W as itself an RCA, is consistent with the postulate of “conscious realism” introduced
 317 by Hoffman and Prakash (2014), who employ RCAs in their discussion of multi-agent com-
 318 binations without introducing this specific terminology. More explicitly, conscious realism
 319 is the ontological claim that the “world” is composed entirely of reduced conscious agents,
 320 and hence can be represented as a network of interacting RCAs as discussed in more detail
 321 in §3.2 below. Conscious realism is effectively, once again, a requirement that the intrinsic
 322 and extrinsic perspectives be mutually consistent: since no RCA can determine that the
 323 internal dynamics of its associated W are non-Markovian from its own intrinsic perspective,
 324 no theoretical, extrinsic-perspective stipulation that its W has non-Markovian dynamics is

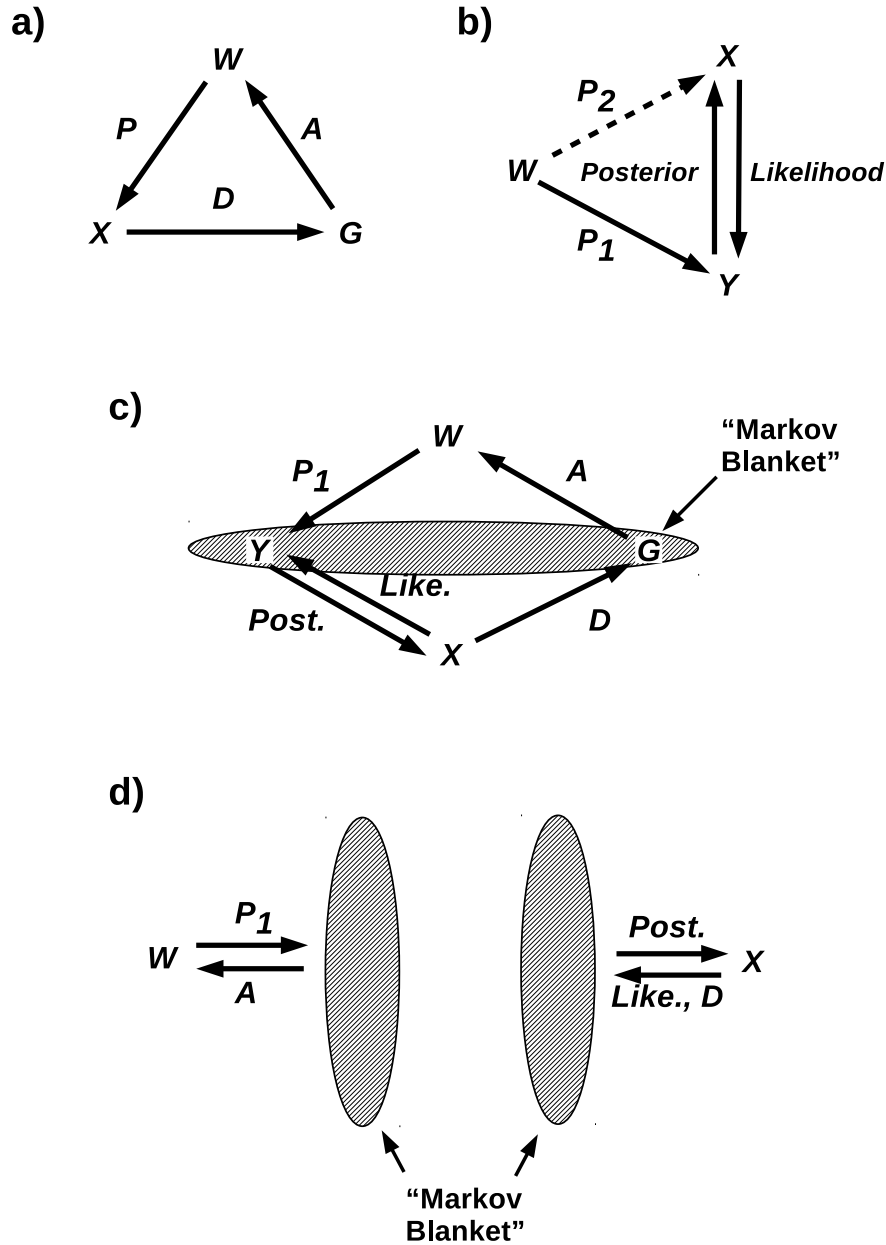
325 allowable. Every occurrence of the symbol W can, therefore, be replaced, as in Fig. 2,
 326 by an RCA. When this is done, all actions - all kernels A - act directly on the experience
 327 spaces X of other RCAs as shown in Fig. 2. If it is possible to consider any arbitrary
 328 system - any directed subgraph comprising sets and kernels - as composing a CA from the
 329 extrinsic perspective, then it is also possible, from the intrinsic perspective of any one of
 330 the RCAs involved, to consider the rest of the network as composing a single RCA with
 331 which it interacts.



332 *Fig. 2:* Representation of an interaction between two RCAs as a labelled di-
 333 rected graph (*cf.* Hoffman and Prakash, 2014, Fig. 2). Note that consistency
 334 requires that the actions A possible to the lower RCA must be the same as the
 335 perceptions P possible for the upper RCA and vice-versa.

336 2.3 Bayesian inference and the Markov blanket

337 As emphasized above, the set X represents the set of possible *experiences* of a conscious
 338 agent within the CA framework. In the case of human beings, including even neonates
 339 (e.g. Rochat, 2012; see also §4 below), such experiences invariably involve *interpretation*
 340 of raw sensory input, e.g. of photoreceptor or hair-cell excitations. It is standard to model
 341 interpretative inferences from raw sensory input or “images” in some modality to expe-
 342 rienced “scene interpretations” (to use visual language) using Bayesian Decision Theory
 343 (BDT; reviewed e.g. by Maloney and Zhang, 2010). In recognition of the fact that such
 344 inferences are executed by the perceiving organism and are hence subject to the constraints
 345 of an evolutionary history, Hoffman and Singh (2012) introduced the framework of Com-
 346 putational Evolutionary Perception (CEP) shown in Fig. 3b. This framework differs from
 347 many formulations of BDT by emphasizing that both posterior probability distributions
 348 and likelihood functions are generated within the organism. The posterior distributions,
 349 in particular, are *not* generated directly by the world W (see also Hoffman, Singh and
 350 Prakash, 2015).



351 *Fig. 3:* Relation between the current CA framework and the “Markov blanket”
 352 formalism of Friston (2013). a) The canonical CA, *cf.* Fig. 1. b) The “Compu-
 353 tational Evolutionary Perception” (CEP) extension of Bayesian decision theory
 354 developed by Hoffman and Singh (2012). Here the set Y is interpreted as a set of
 355 “images” and the set X is interpreted as a set of “scene interpretations,” consis-
 356 tent with the interpretation of X in the CA framework. The map $P_2 : W \mapsto X$

357 is induced by the composition of the “raw” input map P_1 with the posterior-
358 map - likelihood-map loop. c) Identifying P in the CA framework with P_2 in
359 the CEP formalism replaces the canonical CA with a four-node graph. Here the
360 sets Y and G jointly constitute a Markov blanket as defined by Friston (2013).
361 d) Both W and X can be regarded as interacting bi-directionally with just their
362 proximate “surfaces” of the Markov blanket comprising Y and G . The blan-
363 ket thus isolates them from interaction with each other, effectively acting as an
364 interface in the sense defined by ITP.

365 The CEP framework effectively decomposes the kernel P of a CA (Fig. 3a) into the com-
366 position of a mapping P_1 from W to a space Y of “raw” perceptual images with a map
367 (labelled B in Hoffman, Singh and Prakash, 2015, Fig. 4) corresponding to the construc-
368 tion of a posterior probability distribution on X . The state of the image space Y depends,
369 in turn, on the state of X via the feedback of a Bayesian likelihood function; hence the
370 embedded posterior - likelihood loop provides the information exchange between prior and
371 posterior distributions needed to implement Bayesian inference. The Bayesian likelihood
372 serves, in effect, as the perceiving agent’s implicit “model” of the world as it is seen via the
373 image space Y .

374 As shown by Pearl (1988), any set of states that separates two other sets of states from each
375 other in a Bayesian network can be considered a “Markov blanket” between the separated
376 sets of states (*cf.* Friston (2013)). The disjoint union $Y \sqcup G$ of Y and G separates the
377 sets W and X in Fig. 3b in this way; hence $Y \sqcup G$ constitutes a Markov blanket between
378 W and X (*cf.* Friston, 2013, Fig. 1). Each of W and X can be regarded as interacting
379 bidirectionally, via Markov processes, with a “surface” of the Markov blanket, as shown in
380 Fig. 3d. The blanket therefore serves as an “interface” in the sense required by ITP: it
381 provides an indirect representation of W to X that is constructed by processes to which X
382 has no independent access. Consistent with the assumption of conscious realism above, this
383 situation is completely symmetrical: the blanket also provides an indirect representation of
384 X to W that is constructed by processes to which W has no independent access. The role
385 of the Markov blanket in Fig. 3d is, therefore, exactly analogous to the role of the second
386 agent in Fig. 2. The composed Markov kernel $D'A$ in Fig. 2 represents, in this case, the
387 internal dynamics of the blanket.

388 Friston (2013) argues that any random ergodic system comprising two subsystems separated
389 by a Markov blanket can be interpreted as minimizing a variational free energy that can, in
390 turn, be interpreted in Bayesian terms as a measure of expectation violation or “surprise.”
391 This Bayesian interpretation of “inference” through a Markov blanket is fully consistent
392 with the model of perceptual inference provided by the CEP framework. Conscious agents as
393 described here can, therefore, be regarded as free-energy minimizers as described by Friston
394 (2010). This formal as well as interpretational congruence between the CA framework and
395 the free-energy principle (FEP) framework of Friston (2010) is explored further below,
396 particularly in §3.3 and §4.3.

2.4 Effective propagator and master equation

From the intrinsic perspective of a particular CA, experience consists of a sequence of states of X , each of which is followed by an action of D and a “tick” of the internal counter t . The sequence of transitions between successive states of X can be regarded as generated by an effective propagator $T_{\text{eff}} : \mathcal{M}_X(t) \rightarrow \mathcal{M}_X(t+1)$, where $\mathcal{M}_X(t)$ is the collection of probability measures on X at each “time” t defined by the internal counter. This propagator satisfies, by definition, a master equation that, in the discrete t case, is the Chapman-Kolmogorov equation: If μ_t is the probability distribution at time t , then $\mu_{t+1} = T_{\text{eff}}\mu_t$.

The propagator T_{eff} cannot, however, be characterized from the intrinsic perspective: all that is available from the intrinsic perspective is the current state $X(t)$, including, as discussed in §4 below, the current states of any memories contained in $X(t)$. From the extrinsic perspective, the structure of T_{eff} depends on the structure of the world W . Here again, the assumption of conscious realism and hence the ability to represent any W as a second agent as shown in Fig. 2 is critical. In this case, $T_{\text{eff}} = PD'AD$, where in the general case the actions of each of these operators at each t depend on the initial, $t = 0$ state of the network. As discussed above, the P and D kernels within this composition can be regarded as specifying the interaction between X and a Markov blanket with internal dynamics $D'A$. The claim that T_{eff} is a Markov process on X is then just the claim that the composed kernel $PD'AD$ is Markovian, as kernel composition guarantees it must be. As Friston, Levin, Sengupta and Pezzulo (2015) point out, the Markov blanket framework “only make(s) one assumption; namely, that the world can be described as a random dynamical system” (p. 9). Both the above representation of T_{eff} and the Chapman-Kolmogorov equation $\mu_{t+1} = T_{\text{eff}}\mu_t$ are independent of the structure of the Markov blanket, which as discussed in §3.2 below can be expanded into an arbitrarily-complex networks of RCAs, provided this condition is met.

For simplicity, we adopt in what follows the assumption that all relevant Markov kernels, and therefore the propagator T_{eff} , are homogeneous and hence independent of t for any agent under consideration. As discussed further below, this assumption imposes interpretations of both evolution (§3.3) and learning (§4.3) as processes that change the occupation probabilities of states of X and G but do not change any of the kernels P , D or A . This interpretation can be contrasted with that of typical machine learning methods, and in particular, typical artificial neural network methods, in which the outcome of learning is an altered mapping from input to output. The current interpretation is, however, consistent with Friston’s (2010; 2013) characterization of free-energy minimization as a process that maintains homeostasis. In the current framework, the maintenance of homeostasis corresponds to the maintenance of an *experience* of homeostasis, i.e. to continued high probabilities of occupation of particular components of the state of X . Both evolution and learning act to maintain homeostasis and hence maintain these high state-occupation probabilities. This idea that maintenance of homeostasis is signalled by maintaining an experience of homeostasis is consistent with the conceptualization of affective state as an experience-marker of a physiological, and particularly homeostatic state (Damasio, 1999;

439 Peil, 2015). As noted earlier, no assumption that such experiences are reportable by any
440 particular, e.g. verbal behavior are made (see also §3.3, 4.4 below).

441 **3 W from the extrinsic perspective: RCA networks** 442 **and dynamic symmetries**

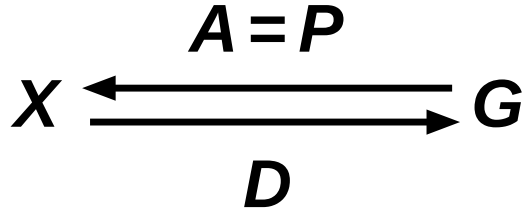
443 **3.1 Symmetric interactions**

444 From the extrinsic perspective, a CA is a syntactic construct comprising three distinct sets
445 of states and three Markovian kernels between them as shown in Fig. 1. We begin here
446 to analyze the behavior of such constructs, starting below with the simplest CA network
447 and then generalizing (§3.2) to networks of arbitrary complexity. Familiar concepts from
448 social-network theory emerge in this setting, and provide (§3.3) a natural characterization
449 of “fitness” for CAs.

450 Here and in what follows, we assume that each of the relevant σ -algebras contains all
451 singleton subsets of its respective underlying set. We call a Markovian kernel “punctual,”
452 i.e. non-dispersive, if the probability measures it assigns are Dirac measures, i.e. measures
453 concentrated on a singleton subset. In this case, P can be regarded as selecting a single
454 element x from X , and can therefore be identified with a *function* from $W \times X$ to X .
455 The punctual kernels between any pair of sets are the extremal elements of the set of
456 all kernels between those sets provided the relevant σ -algebras contain all of the singleton
457 subsets as assumed above; hence characterizing their behavior in the discrete case implicitly
458 characterizes the behavior of all kernels in the set. The punctual kernels of a network of
459 interacting RCAs specify, in particular, the extremal dynamics of the network. Conscious
460 realism entails the purely syntactic claim that the graphs shown in Figs. 1 and 2 are
461 interchangeable as discussed above; the world W can, therefore, be regarded as an arbitrarily-
462 complex network of interacting RCAs, subject only to the constraint that the A and P
463 kernels of the interacting RCAs can be identified (Hoffman and Prakash, 2014).

464 The simplest CA network is a dyad in which $W = X \sqcup G$, where as above the notation
465 $X \sqcup G$ indicates the disjoint union of X with G , and $A = P$; it is shown in Fig. 4. This
466 dyad acts on its own X ; its perceptions are its actions. From a purely formal perspective,
467 this dyad is isomorphic to the $X - Y$ dyad of the CEP framework (Fig. 3b); it is also
468 isomorphic to the interaction of X with its proximal “surface” of a Markov blanket separ-
469 ating it from W (Fig. 3d). Investigating the behavior of this network over time requires
470 specifying, from the extrinsic perspective, the state spaces and operators. The simplest
471 case is the *symmetric interaction* in which the two state spaces are identical. If both X and
472 G are taken to contain just one bit, the four possible states of the network can be written
473 as $|00\rangle, |01\rangle, |10\rangle$ and $|11\rangle$. Here we will represent these states by the orthogonal (column)
474 vectors $(1, 0, 0, 0)^T, (0, 1, 0, 0)^T, (0, 0, 1, 0)^T$ and $(0, 0, 0, 1)^T$, respectively. The simplest ker-
475 nels $D : X \times G \rightarrow G$ and $A : G \times X \rightarrow X$ are punctual. Let $x(t)$ and $g(t)$ denote the

476 state of X and G , respectively, at time t . We slightly abuse the notation and use the letter
 477 D to refer to the operator $I_X \otimes D : X(t) \times G(t) \rightarrow X(t+1) \times G(t+1)$, where I_X is the
 478 Identity operator on X . This D leaves the state x of X unchanged but changes the state
 479 of G to $g(t+1) = D(x(t), g(t))$. Similarly, we will use the letter A to refer to the operator
 480 $A \otimes I_G : X(t) \times G(t) \rightarrow X(t+1) \times G(t+1)$, where I_G is the identity operator on G . This
 481 A leaves the state g of G unchanged, but changes the state of X to $x(t+1) = A(g(t), x(t))$.
 482 Note that in this representation, D and A are both executed each time the “clock ticks.”



483 *Fig. 4:* The simplest possible CA network, the dyad in which $W = X \sqcup G$.

484 To reiterate, the decision operator D acts on the state of G but leaves the state of X
 485 unchanged, i.e. $X(t+1) = X(t)$. Only four Markovian operators with this behavior exist.
 486 These are the identity operator,

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix};$$

487 the NOT operator,

$$\mathbf{N}_D = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix};$$

488 the controlled-NOT (cNOT) operator that flips the G bit when the X bit is 0,

$$\mathbf{C}_{D0} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix};$$

489 and the cNOT operator that flips the G bit when the X bit is 1,

$$\mathbf{C}_{D1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

490 The action operator A acts on the state of X but leaves the state of G unchanged, i.e.
 491 $G(t+1) = G(t)$. Again, only four Markovian operators with this behavior exist. These are
 492 the identity operator \mathbf{I} defined above, the NOT operator,

$$\mathbf{N}_A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix};$$

493 the cNOT operator that flips the X bit when the G bit is 0,

$$\mathbf{C}_{A0} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix};$$

494 and the cNOT operator that flips the X bit when the G bit is 1,

$$\mathbf{C}_{A1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

495 In principle, distinct CAs with single-bit X and G could be constructed with any one of
 496 the four possible D operators and any one of the four possible A operators. The CA in
 497 which both operators are identities is trivial: it never changes state. The CA in which both
 498 operators are NOT operators is the familiar bistable multivibrator or “flip-flop” circuit. It
 499 is also interesting, however, to consider the abstract entity – referred to as a “participator”
 500 in Bennett, Hoffman and Prakash (1989) – in which X and G are fixed at one bit and all
 501 possible D and A operators can be employed. The dynamics of this entity are generated by
 502 the operator compositions DA and AD . There are 24 distinct compositions of the above
 503 7 operators, which form the Symmetric Group on 4 objects, S_4 . This group appears in a
 504 number of geometric contexts and is well characterized; the CA dynamics with this group of
 505 transition operators include limit cycles, i.e. cycles that repeatedly revisit the same states,
 506 of lengths 1 (the identity operator \mathbf{I}), 2, 3 and 4. Hence there are 24 distinct CAs having
 507 the form of Fig. 3 but with different choices for D and A , with behavior ranging from
 508 constant ($D = A = \mathbf{I}$) to limit cycles of length 4.

509 It is important to emphasize that there is no sense in which the 1-bit dyad *experiences* the
510 potential complexity of its dynamics, or in which the experience of a 1-bit dyad with one
511 choice of D and A operators is any different from the experience of a 1-bit dyad with another
512 choice of operators. Any 1-bit dyad has only two possible experiences, those tokened by $|0\rangle$
513 and $|1\rangle$. The addition of memory to a CA in order to enable it to experience a *history* of
514 states and hence relations between states from its own intrinsic perspective is discussed in
515 §4 below.

516 The Identity and NOT operators can be expressed as “forgetful” kernels, i.e. kernels that
517 do not depend on the state at t of their target spaces, $D : X(t) \rightarrow G(t + 1)$ and $A : G(t) \rightarrow X(t + 1)$ but the cNOT operators cannot be; hence the forgetful kernels introduced
518 by Hoffman and Prakash (2014) have less representational power than the state-dependent
519 kernels employed in the current definition of a CA. It is also worth noting that the standard
520 AND operator taking $x(t)$ and $g(t)$ to $x(t + 1) = x(t)$ and $g(t + 1) = x(t)$ AND $g(t)$ may
521 be represented as:
522

$$\mathbf{AND}_{\mathbf{G}} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

523 and the corresponding OR operator taking $x(t)$ and $g(t)$ to $x(t+1) = x(t)$ and $g(t+1) = x(t)$
524 OR $g(t)$ may be represented as:

$$\mathbf{OR}_{\mathbf{G}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}.$$

525 The value of $G(t)$ cannot be recovered following the action of either of these operators; they
526 are therefore logically irreversible. As each of the matrix representations of these operators
527 has a row of all zeros, they are not Markovian. The logically irreversible, non-Markovian
528 nature of these operators has, indeed, been a primary basis of criticisms of artificial neural
529 network and dynamical-system models of cognition; Fodor and Pylyshyn (1988), for
530 example, criticize such models as unable, in principle, to replicate the compositionality of
531 Boolean operations in domains such as natural language. The standard AND operator
532 can, however, be implemented reversibly by adding a single ancillary z bit to X , fixing its
533 value at 0, and employing the Toffoli gate that maps $[x, y, z]$ to $[x, y, (x \text{ AND } y) \text{ XOR } z]$,
534 where XOR is the standard exclusive OR (Toffoli, 1980). The Toffoli gate preserves the
535 values of x and y and allows the value of z to be computed from the values of x and y ;
536 hence it is reversible and can, therefore, be represented as a punctual Markovian kernel.
537 The standard XOR operator employed in the Toffoli gate is equivalent to a cNOT. As any
538 universal computing formalism must be able to compute AND, the 1-bit dynamics of Fig.
539 4 is not computationally universal. The Toffoli gate is, however, computationally universal,

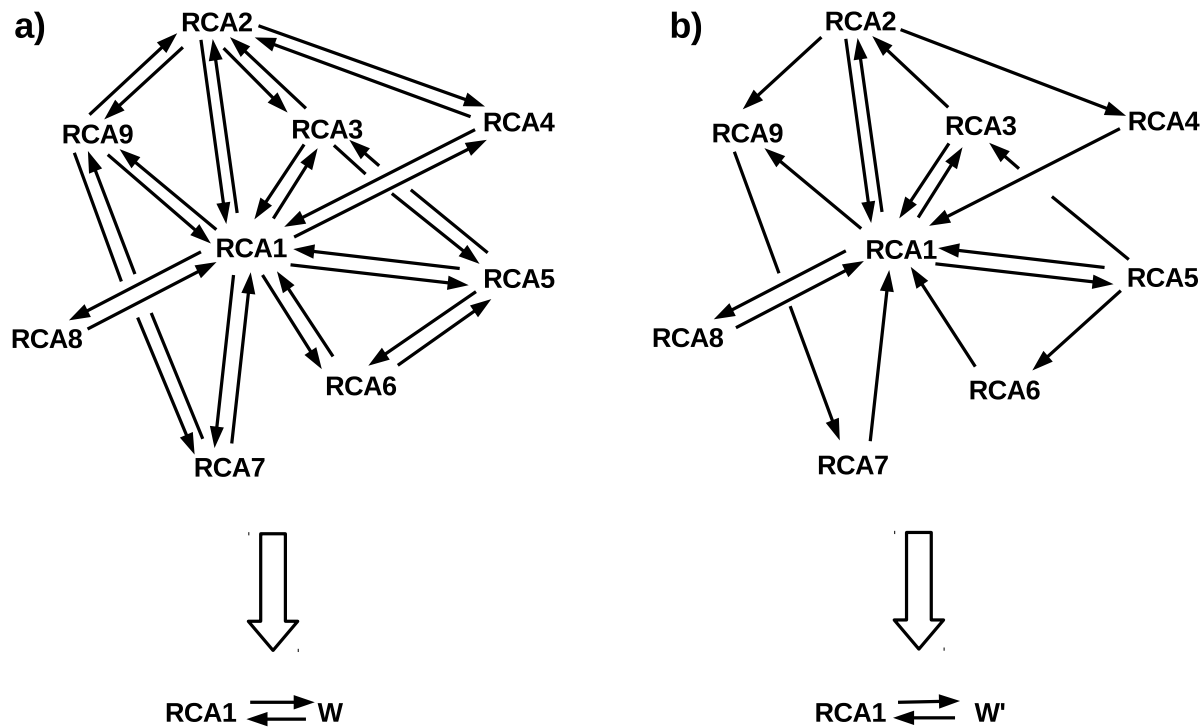
540 so adding a single ancillary bit set to 0 to each space in Fig. 4 is sufficient to achieve
541 universality.

542 Two distinct graphs representing symmetric, punctual CA interactions have 4 bits in total
543 and hence 16 states: the graph shown in Fig. 2 where each of X, G, X' and G' contains one
544 bit and the graph shown in Fig. 4 in which each of X and G contains 2 bits. These graphs
545 differ from the intrinsic as well as the extrinsic perspectives: in the former case each agent
546 experiences only $|0\rangle$ or $|1\rangle$ – i.e. has the same experience as the 1-bit dyad – while in the
547 latter case the agent has the richer experience $|00\rangle, |01\rangle, |10\rangle$ or $|11\rangle$. The dynamics of the
548 participator with the first of these structures has been exhaustively analyzed; it has the
549 structure of the affine group $AGL(4,2)$. Further analyses of the dynamics of these simple
550 systems, including explicit consideration of the behavior of the t counters, is currently
551 underway and will be reported elsewhere.

552 While the restriction to punctual kernels simplifies analysis, systems in which perception,
553 decision and action are characterized by dispersion will have non-punctual kernels P, D and
554 A . It is worth noting that from the extrinsic, theorist’s perspective, such dispersion exists
555 by stipulation: the kernels P, D and A characterizing a particular CA within a particular
556 situation being modelled are stipulated to be stochastic. The probability distributions on
557 states of X, G and W that they generate are, from the theorist’s perspective, distributions
558 of objective probabilities: they are stipulated “from the outside” as fixed components of the
559 theoretical model. As will be discussed in §4 below, these become *subjective* probabilities
560 when viewed from the intrinsic perspective of any observer represented within such a model.
561 However as noted earlier, ITP forbids any CA from having observational access to its own
562 P, D , or A kernels; hence no CA can determine by observation that its kernels are non-
563 punctual.

564 3.2 Asymmetric interactions and RCA combinations

565 While symmetric interactions are of formal interest, a “world” containing only two sub-
566 systems of equal size has little relevance to either biology or psychology. Real organisms
567 inhabit environments much larger and richer than they are, and are surrounded by other
568 organisms of comparable size and complexity. The realistic case, and the one of interest
569 from the standpoint of ITP, is that in which the σ -algebra \mathcal{W} is much finer than either
570 \mathcal{X} or \mathcal{G} . This asymmetrical interaction can be considered effectively bandwidth-limited by
571 the relatively small encoding capacities of \mathcal{X} and \mathcal{G} . Representing the two-RCA interaction
572 shown in Fig. 2 by the shorthand notation $RCA1 \leftrightarrow RCA2$, this more realistic situation can
573 be represented as in Fig. 5, in which no assumptions are made about the relative “sizes”
574 of the RCAs or the dimensionality of the Markovian kernels involved.

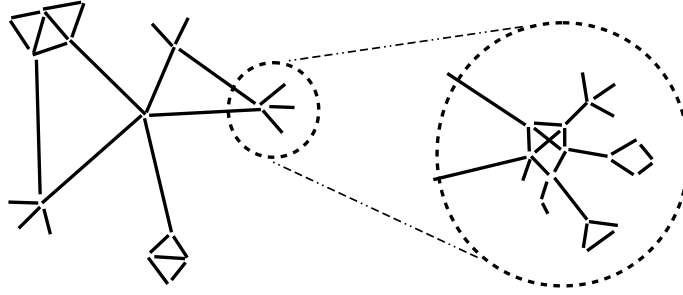


575 *Fig. 5:* a) Nine bidirectionally interacting RCAs, equivalent to a single RCA
 576 interacting with its “world” W and hence to a single CA. b) A network similar
 577 to that in a), except that some interactions are not bidirectional. Here again,
 578 the RCA network is equivalent to a single RCA interacting with a structurally
 579 distinct “world” W’ and hence to a distinct single CA. In general, RCA networks
 580 of either kind are asymmetric for every RCA involved.

581 When applied to the multi-RCA interaction in Fig. 5, consistency between intrinsic and
 582 extrinsic perspectives requires that when a theorist’s attention is focussed on any single
 583 RCA, the other RCAs together can be considered to be the “world.” If attention is focussed
 584 on RCA1, for example, it must be possible to regard the subgraph comprising RCA2 - RCA9
 585 as the “world” W (Fig. 5a) and the entire network as specifying a single CA in the canonical
 586 form of Fig. 1. As every RCA interacts bidirectionally with its “world,” any directed path
 587 within an RCA network must be contained within a closed directed path. These paths
 588 do not, however, all have to be bidirectional; the RCA network in Fig. 5b can equally
 589 well be represented in the canonical form of Fig. 1. The “worlds” of Fig. 5a and Fig.
 590 5b have distinct structures from the extrinsic perspective. However, ITP requires that the
 591 interaction between RCA1 and its “world” does not determine the internal structure of

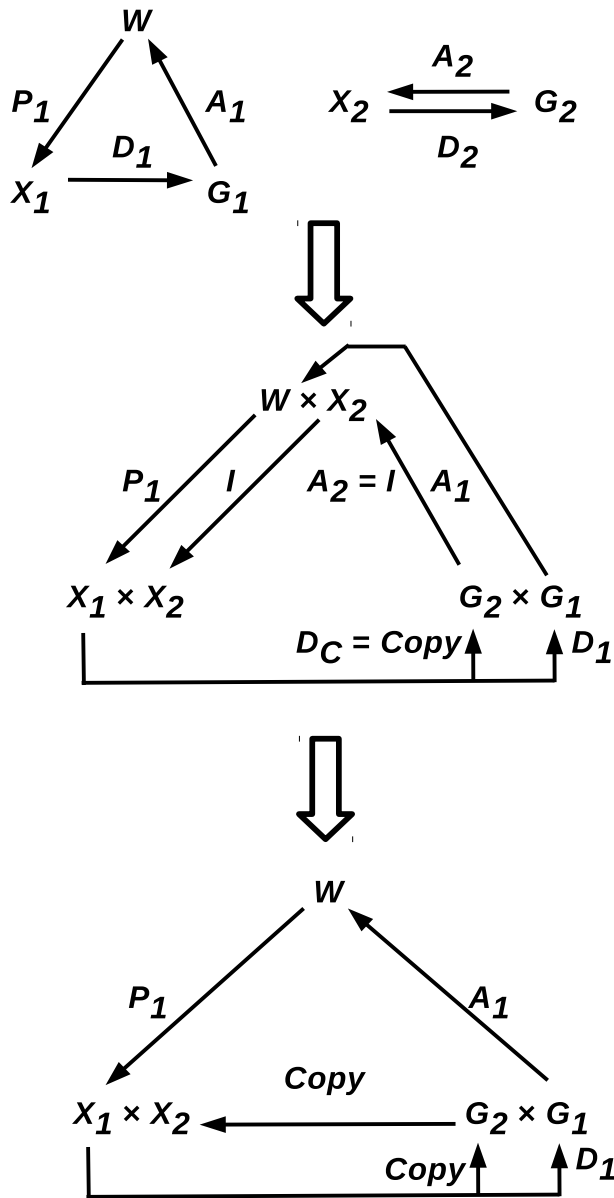
592 the “world”; indeed an arbitrarily large number of alternative structures could produce
593 the same inputs to RCA1 and hence the same sequence of experiences for RCA1. RCA1
594 cannot, in particular, determine what other RCA(s) it is interacting with at any particular
595 “time” t as measured by its counter, or determine whether the structure or composition of
596 the network of RCAs with which it is interacting changes from one value of t to the next.
597 This lack of transparency renders the “world” of any RCA a “black box” as defined by
598 classical cybernetics (Ashby, 1956): a system with an internal structure under-determined,
599 in principle, by finite observations. Even a “good regulator” (Conant and Ashby, 1970) can
600 only regulate a black box to the extent that the behavior of the box remains within the
601 bounds for which the regulator was designed; whether a given black box will do so is always
602 unpredictable even in principle. From the intrinsic perspective of the “world,” the same
603 reasoning renders RCA1 a black box; hence consistency between perspectives requires that
604 any RCA - and hence any CA - for which the sets X and G are not explicitly specified be
605 regarded as potentially having an arbitrarily rich internal structure.

606 In general, consistency between intrinsic and extrinsic perspectives requires that any ar-
607 bitrary connected network of RCAs can be considered to be a single canonical-form CA;
608 for each RCA in the network, all of the other RCAs in the network, regardless of how
609 they are connected, together form of “world” of that RCA. Non-overlapping boundaries
610 can, therefore, be drawn arbitrarily in a network of interacting RCAs and the RCAs within
611 each of the boundaries “combined” to form a smaller network of interacting RCAs, with
612 a single canonical-form CA or $X - G$ dyad as the limiting case in which all RCAs in
613 the network have been combined. Connected networks that characterize gene regulation
614 (Agrawal, 2002), protein interactions (Barabási and Oltvai, 2004), neurocognitive archi-
615 tecture (Bassett and Bullmore, 2006), academic collaborations (Newman, 2001) and many
616 other phenomena exhibit dynamic patterns including preferential attachment (new connec-
617 tions are preferentially added to already well-connected nodes; Barabási and Albert, 1999)
618 and the emergence of small-world structure (short minimal path lengths between nodes
619 and high clustering; Watts and Strogatz, 1998). Such networks typically exhibit “rich
620 club” connectivity, in which the most well-connected nodes at one scale form a small-world
621 network at the next-larger scale (Colizza, Flammini, Serrano and Vespignani, 2006); the
622 human connectome provides a well-characterized example (van den Heuvel and Sporns,
623 2011). Networks in which connectivity structure is, on average, independent of scale are
624 called “scale-free” (Barabási, 2009); such networks have the same structure, on average,
625 “all the way down.” As illustrated in Fig. 6, scale-free structures approximate hierarchies;
626 “zooming in” to a node in a small-world or rich-club network typically reveals small-world
627 or rich-club structure within the node. However, these networks allow the “horizontal”
628 within-scale connections that a strict hierarchical organization would forbid. Given the
629 prominence of scale-free small-world or rich-club organization in Nature, it is reasonable to
630 ask whether RCA networks can exhibit such structure. In particular, it is reasonable to ask
631 whether interactions between “simple” RCAs can lead to the emergence of more complex
632 RCAs that interact among themselves in an approximately-hierarchical, rich-club network.
633 We consider this question in one particular case in §4 below.



634 *Fig. 6:* “Zooming in” to a node in a rich-club network typically reveals addi-
 635 tional small-world structure at smaller scales. Here the notation has been further
 636 simplified by eliding nodes altogether and only showing their connections.

637 Replication followed by functional diversification ubiquitously increases local complexity in
 638 biological and social systems; processes ranging from gene duplication through organismal
 639 reproduction to the proliferation of divisions in corporate organizations exhibit this process.
 640 The simplest case, for an RCA, is to replicate part or all of the experience set X ; as
 641 will be shown below (§4.2), this operation is the key to building RCAs with memory.
 642 Let $[(X_1, \mathcal{X}_1), (G_1, \mathcal{G}_1), D_1, t_1]$ be an RCA interacting with W via A_1 and P_1 kernels. Let
 643 $[(X_2, \mathcal{X}_2), (G_2, \mathcal{G}_2), D_2, A_2, t_2]$ be a dyad as shown in Fig. 4. Setting $t_1 = t_2 = t$, a new
 644 RCA whose “world” is the Cartesian product $W \times X_2$ can be constructed by taking the
 645 Cartesian products of the sets X_1 and X_2 and G_1 and G_2 respectively, as illustrated in
 646 Fig. 7, and defining product σ -algebras of \mathcal{X}_1 and \mathcal{X}_2 and \mathcal{G}_1 and \mathcal{G}_2 respectively. If all the
 647 kernels are left fixed, these product operations change nothing; they merely put the the
 648 original RCA and the dyad “side by side” in the new, combined RCA. We can, however,
 649 create an RCA with qualitatively new behavior by redefining one or more of the kernels;
 650 the “combination” process in this case significantly alters the behavior of one or both of the
 651 RCAs being “combined.” For example, we can specify a new punctual kernel D'_2 that acts
 652 on the X_1 component instead of the X_2 component of $X_1 \times X_2$, i.e. $D'_2 : X_1 \rightarrow G_2$. Consider,
 653 for example, the RCA that results if D_2 is replaced by a kernel $D'_2 = D_C$ that simply *copies*,
 654 at each t , the current value x_1 of X_1 to G_2 . If the kernel A_2 is set to the Identity I , the
 655 value x_1 will be copied, by A_2 , back to X_2 on each cycle, as shown in Fig. 7. In this case,
 656 the experience of the “combined” CA at each t has two components: the current value of
 657 x_1 and the previous value of x_1 , now “stored” as the value x_2 . This “copying” construction
 658 will be used repeatedly in §4 below to construct agents with progressively more complex
 659 memories. Note that for these memories to be *useful* in the sense of affecting choices of
 660 action, the kernel D_1 must be replaced by one that also depends on the “memory” X_2 .



661 *Fig. 7:* A CA as shown in Fig. 1 and a dyad as shown in Fig. 3 can be
 662 “combined” to form a composite CA with a simple, one time-step short-term
 663 memory by replacing the decision kernel D_2 of the dyad with a kernel D_C that
 664 “copies” the state $x_1(t)$ to $g_2(t+1)$ and setting the action kernel A_2 of the dyad
 665 to the Identity I . The notation can be simplified by eliding the explicit $W \times X_2$
 666 to W and treating the I^2 operation on G_2 as a feedback operation “internal to”
 667 the RCA, as shown in the lower part of the figure. Note that the composite

668 CA produced by this “combination” process has qualitatively different behavior
669 than either of the CAs that were combined to produce it.

670 The construction shown in Fig. 7 suggests a general feature of RCA networks: asymmetric
671 kernels characterize the interactions between typical RCAs and W , but also characterize
672 “internal” interactions that give RCAs additional structure. Such kernels may lose infor-
673 mation and hence “coarse-grain” experience. If RCA networks are indeed scale-free, one
674 would expect asymmetric interactions to be the norm: wherever the RCA-of-interest to W
675 boundary is drawn, the networks on both sides of the boundary would have asymmetric
676 kernels and complex internal organization. If this is the case, the notion of combining ex-
677 perience qualia underlying classic statements of the “combination problem” by William
678 James, Thomas Nagel and many others (for review, see Hoffman and Prakash, 2014) appears
679 too limited. There is no reason, in general, to expect “lower-level” experiences to combine
680 into “higher-level” experiences by Cartesian products. An initially diffuse, geometry-less
681 experience of “red” and an initially color-less experience of “circle,” for example, can be
682 combined to an experience of “red circle” only if the combination process forces the diffuse
683 redness into the boundary defined by the circle. This is not a mere Cartesian product; the
684 redness and the circularity are not merely overlaid or placed next to each other. While
685 Cartesian products of experiences allow recovery of the individual component experiences
686 intact; arbitrary operations on experiences do not. The “combination” operations of inter-
687 est here instead introduce scale-dependent constraints of the type Polanyi (1968) shows are
688 ubiquitous in biological systems (*cf.* Rosen, 1986; Pattee, 2001). Such constraints introduce
689 qualitative novelty. Once the redness has been forced into the circular boundary, for exam-
690 ple, its original diffuseness is not recoverable: the red circle is a qualitatively new construct.
691 Asymmetric kernels, in general, render higher-level agents and their higher-level experiences
692 irreducible. Human beings, for example, experience edges and faces, but early-visual edge
693 detectors do not experience edges and “face detectors” in the Fusiform Face Area do not
694 experience faces. von Uexküll (1957), Gibson (1979) and the embodied cognition movement
695 have made this point previously; the present considerations provide a formal basis for it
696 within the theoretical framework of ITP.

697 3.3 Connectivity and fitness

698 As noted in the Introduction, ITP was originally motivated by evolutionary game simula-
699 tions showing that model organisms with perceptual systems sensitive only to fitness drove
700 model organisms with veridical perceptual systems to extinction (Mark, Marion and Hoff-
701 man, 2010). In these simulations, “fitness” was an arbitrarily-imposed function dependent
702 on the states of both the model environment and the model organism. The assumption of
703 conscious realism, however, requires that it be possible to regard the environment of any
704 organism, i.e. of any agent, as itself an agent and hence itself subject to a fitness function.
705 From a biological perspective, this is not an unreasonable requirement: the environments of
706 all organisms are populated by other organisms, and organism - organism interactions, e.g.

707 predator - prey or host - pathogen interactions, are key determiners of fitness. In the case
708 of human beings, the hypothesis that interactions with conspecifics are the *primary* de-
709 terminant of fitness motivates the broadly-explanatory “social brain hypothesis” (Adolphs,
710 2003, 2009; Dunbar, 2003; Dunbar and Shultz, 2007) and much of the field of evolutionary
711 psychology. If interactions between agents determine fitness, however, it should be possible
712 to derive a representation of fitness entirely *within* the CA formalism. As the minimiza-
713 tion of variational free energy or Bayesian surprise has a natural interpretation in terms of
714 maintenance of homeostasis (Friston, 2013; Friston, Levin, Sengupta and Pezzulo, 2015),
715 the congruence between the CA and FEP frameworks discussed above also suggests that
716 a fully-internal definition of fitness should be possible. Here we show that an intuitively-
717 reasonable definition of fitness not only emerges naturally within the CA framework, but
718 also corresponds to well-established notions of centrality in complex networks.

719 The time parameter t characterizing a CA is, as noted earlier, not an “objective” time but
720 rather an observer-specific, i.e. CA-specific time. The value of t is, therefore, intimately
721 related to the fitness of the CA that it characterizes: a CA with a small value of t has not
722 survived, i.e. not maintained homeostasis for very long by its own internal measure, while
723 a CA with a large value of t has survived a long time. Hence it is reasonable to regard
724 the value of t as a *prima facie* measure of fitness. As t is internal to the CA, this measure
725 is internal to the CA framework. It is, however, not in general an *intrinsic* measure of
726 fitness, as CAs in general do not include an explicit representation of the value of t within
727 the experience space X . From a formal standpoint, t measures the number of executions
728 of D . As D by definition executes whenever a new experience is received into X , the value
729 of t effectively measures the *number of inputs* that a CA has received. To the extent that
730 D selects non-null actions, the value of t also measures the number of outputs that a CA
731 generates.

732 From the intrinsic perspective, a particular RCA cannot identify the source of any particular
733 input as discussed above; inputs can equivalently be attributed to one single W or to
734 a collection of distinct other RCAs, one for each input. The value of t can, therefore,
735 without loss of generality be regarded as measuring the *number of input connections* to
736 other RCAs that an given RCA has. The same is clearly true for outputs: from the
737 intrinsic perspective, each output may be passed to a distinct RCA, so t provides an upper
738 bound on output connectivity. From the extrinsic perspective, the connectivity of any RCA
739 network can be characterized; in this case the number of inputs or outputs passed along
740 a directed connection can be considered a “connection strength” label. The value of t
741 then corresponds to the sum of input connection strengths and bounds the sum of output
742 connection strengths.

743 We propose, therefore, that the “fitness” of an RCA within a fixed RCA network can
744 simply be identified with its input connectivity viewed quantitatively, i.e. as a sum of
745 connection-strength labels, from the extrinsic perspective. In this case, a new connection
746 preserves homeostasis to the extent that it enables or facilitates future connections. A
747 new connection that inhibits future connectivity, in contrast, disrupts homeostasis. In
748 the limit, an RCA that ceases to interact altogether is “dead.” If the behavior of the

749 network is monitored over an extrinsic time parameter (e.g. a parameter that counts the
750 total number of messages passed in the network), an RCA that stops sending or receiving
751 messages is dead. The “fittest” RCAs are, in contrast, those that continue to send and
752 receive messages, i.e. those that continue to interact with their neighbors, over the longest
753 extrinsically-measured times. Among these, those RCAs that exchange messages at the
754 highest frequencies for the longest are the most fit.

755 For simple graphs, i.e. graphs with at most one edge between each pair of nodes, the
756 “degree” of a node is the number of incident edges; the input and output degrees are the
757 number of incoming and outgoing edges in a digraph (e.g. Diestel, 2010 or for specific
758 applications to network theory, Börner, Sanyal and Vespignani, 2007). A node is “degree
759 central” or has maximal “degree centrality” within a graph if it has the largest degree;
760 nodes of lower degree have lower degree centrality. These notions can clearly be extended
761 to labelled digraphs in which the labels indicate connection strength; here “degree” becomes
762 the sum of connection strengths and a node is “degree central” if it has the highest total
763 connection strength. Applying these notions to RCA networks with the above definition of
764 fitness, the fitness of an RCA scales with its input degree, and hence with its input degree
765 centrality. Note that a small number of high-strength connections can confer higher degree
766 centrality and hence higher fitness than a large number of low-strength connections with
767 these definitions.

768 In an initially-random network that evolves subject to preferential attachment (Barabási
769 and Albert, 1999), the connectivity of a node tends to increase in proportion to its existing
770 connectivity; hence “the rich get richer” (the “Matthew Effect”; see Merton, 1968). As
771 noted above, this drives the emergence of small-world structure, with the nodes with high-
772 est total connectivity forming a “rich club” with high mutual connectivity. Nodes within
773 the rich club clearly have high degree centrality; they also have high betweenness centrality,
774 i.e. paths between non-rich nodes tend to traverse them (Colizza, Flammini, Serrano and
775 Vespignani, 2006). The identification of connectivity with fitness is obviously quite natu-
776 ral in this setting; the negative fitness consequences of isolation are correspondingly well
777 documented (e.g. Steptoe, Shankar, Demakakos and Wardle, 2013).

778 The identification of fitness with connectivity provides a straightforward solution to the
779 “dark room” problem faced by uncertainty-minimization systems (e.g. Friston, Thornton
780 and Clark, 2012). Dark rooms do not contain opportunities to create or maintain connec-
781 tions; therefore fitness-optimizing systems can be expected to avoid them. This solution
782 complements that of Friston, Thornton and Clark (2012), who emphasize the costs to
783 homeostasis of remaining in a dark room. Here again, interactivity and maintenance of
784 homeostasis are closely coupled.

4 W from the intrinsic perspective: Prediction and effective action

4.1 How can non-veridical perceptions be useful?

The fundamental question posed by the ITP is that of how non-veridical perceptions can be informative and hence *useful* to an organism. As noted in the Introduction, veridical perception is commonly regarded as “absolutely essential” for utility; non-veridical perceptions are considered to be illusions or errors (e.g. Pizlo, Sawada and Steinman, 2014). We show in this section that CAs that altogether lack veridical perception can nonetheless exhibit complex adaptive behavior, an outcome that is once again consonant with that obtained within the free-energy framework (Friston, 2010; 2013). We show, moreover, that constructing a CA capable of useful perception and action in a complex environment leads to predictions about both the organization of long-term memory and the structure of object representations that accord well with observations.

For any particular RCA, the dynamical symmetries described in §3.1 are manifested by repeating patterns of states of X . The question of utility can, therefore, be formulated from the intrinsic perspective as the question of how an RCA can detect, and make decisions based on, repeating patterns of states of its own X . As the complexities of both the agent and the world increase, moreover, the probability of a complete experience - a full state of X - being repeated rapidly approaches zero. For agents such as human beings living in a human-like world, only particular aspects of experience are repeated. Such agents are faced with familiar problems, including perceptual figure-ground distinction, the inference of object persistence and hence object identity over time, correct categorization of objects and events, and context dependence (“contextuality” in the quantum theory and general systems literature; see e.g. Kitto, 2014). Our goal in this section is to show that the CA formalism provides a useful representation for investigating these and related questions. We show, in particular, that the limited syntax of the CA formalism is sufficient to implement memory, predictive coding, active inference, attention, categorization and planning. These functions emerge naturally, moreover, from asking what structure an RCA must have in order for its perceptions to be useful for guiding action within the constraints imposed by ITP. We emphasize that by “useful” we mean useful to the RCA from its own intrinsic perspective, e.g. useful as a guide to actions that lead to experiences that match its prior expectations (*cf.* Friston, 2010).

We explicitly assume that the experiences of any RCA are determinate or “classical”: an RCA experiences just one state of X at each t . From the intrinsic perspective of the RCA, therefore, P is always *apparently* punctual regardless of its extrinsic-perspective statistical structure; from the intrinsic perspective, P specifies what the RCA *does* experience, not just what it *could* experience. The RCA selects, moreover, just one action to take at each t ; hence D is *effectively* punctual, specifying what the RCA does do as opposed to merely what it could do, from the intrinsic perspective. This effective or apparent resolution of a probability distribution into a single chosen or experienced outcome is referred to as the

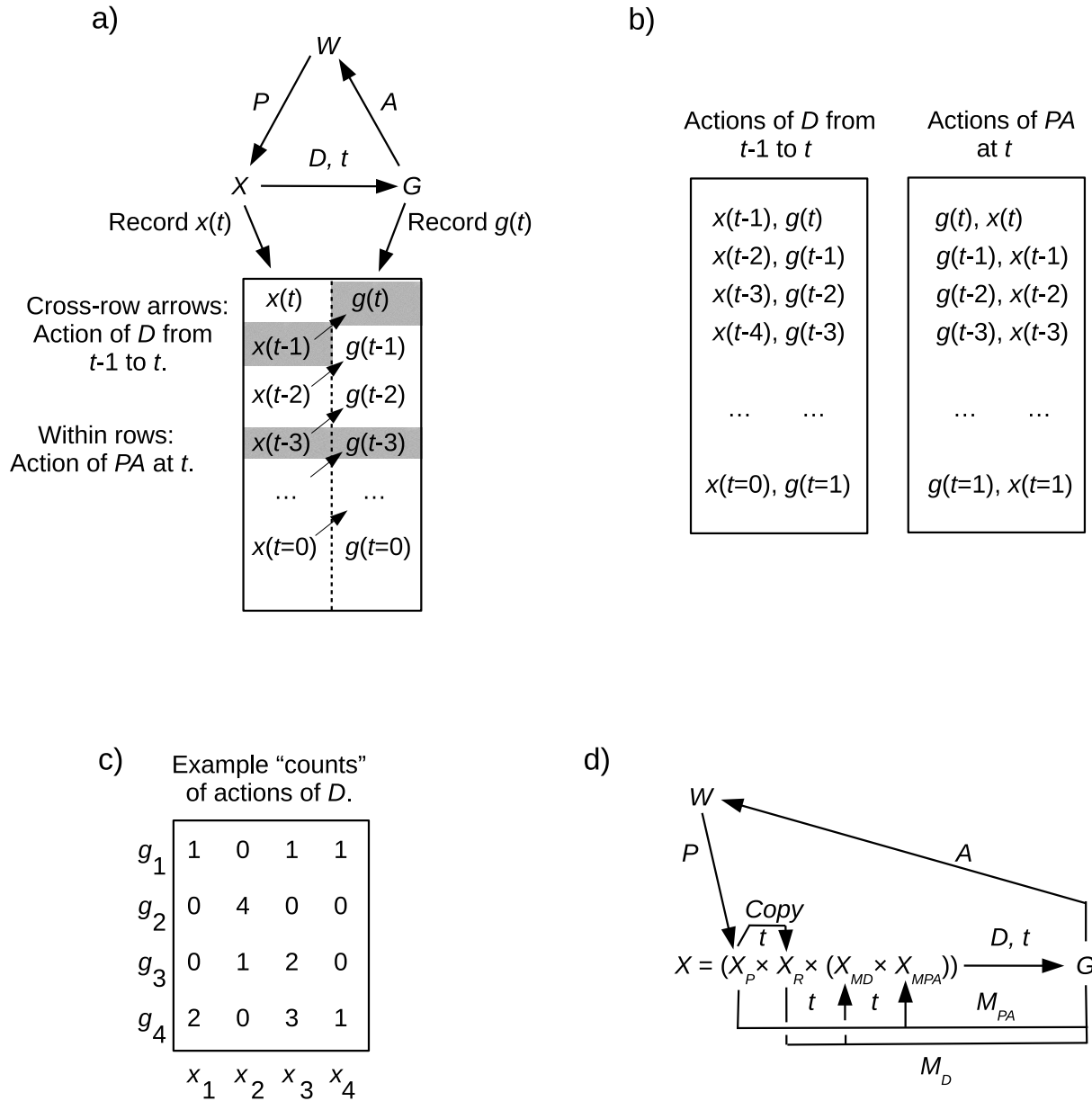
825 “collapse of the wavefunction” in quantum theory (for an accessible and thorough review,
826 see Landsman, 2007) and is often associated with the operation of free will (reviewed by
827 Fields, 2013a). We adopt this association of “collapse” with free will here: the RCA renders
828 P punctual by choosing which of the possibilities offered by W to experience, and renders
829 D punctual by choosing what to do in response. As is the case in quantum theory (Conway
830 and Kochen, 2006), consistency between intrinsic and extrinsic perspectives requires that
831 free will also be attributed to W ; hence we regard W , as an RCA, choosing how to respond
832 to each action A taken by any RCA embedded in or interacting with it. All such choices
833 are regarded as instantaneous. Consistency between internal and external perspectives
834 requires, moreover, that all such choices are unpredictable in principle. An RCA with
835 sufficient cognitive capabilities can, in particular, predict what it *would choose*, given its
836 current state, to do in a particular circumstance, but cannot predict what it *will* do, i.e.
837 what choice it will actually make, when that circumstance actually arises. This restriction
838 on predictions is consonant with a recent demonstration that predicting an action requires,
839 in general, greater computational resources than taking the action (Lloyd, 2012).

840 4.2 Memory

841 Repeating patterns of perceptions are only useful if they can be detected, learned from, and
842 employed to influence action. Within the CA framework, “detecting” something involves
843 awareness of that something; detecting something is therefore a state change in X . Noticing
844 that a current perception repeats a past one, either wholly or in part, requires a memory
845 of past perceptions and a means of comparing the current perception to remembered past
846 perceptions. Both current and past perceptions are states in X , so it is natural to view
847 their comparison as an operation on X . Using patterns of repeated perceptions to influence
848 action requires, in turn, a representation of how perception affects action: an accessible,
849 internal “model” of the D kernel. Consider, for example, an agent with a 1-bit X that
850 experiences only “hungry” and “not hungry” and implements the simple operator, “eat if
851 but only if hungry” as D . This agent has no representation, in X , of the action “eat”; hence
852 it cannot associate hunger with eating, or eating with the relief of hunger. It has, in fact, no
853 representation of any action at all, and therefore no knowledge that it has ever acted. There
854 is no sense in which this agent can learn anything, from its own intrinsic perspective, about
855 W or about its relationship to W . Learning about its relationship to the world requires, at
856 minimum, an ability to experience its own actions, i.e. a representation of those actions in
857 X . This is not possible if X has only one bit.

858 The construction of a memory associating actions with their immediately-following per-
859 ceptions is shown in Fig. 8a. Here as before, t increments when D executes. Note that
860 while each within-row pairing $(g(t), x(t))$ provides a sample and hence a partial model of
861 W ’s response to the choice of $g(t)$, i.e. of the action of the composite kernel PA , each
862 cross-row pairing $(g(t), x(t - 1))$ provides a sample and hence a partial model of the action
863 of D . As noted earlier, no specific assumption about the units of t is made within the CA
864 framework; hence the scope and complexity of the action - perception associations recorded

865 by this memory is determined entirely by the definition, within a particular model, of the
 866 decision kernel D .



867 *Fig. 8:* Constructing a memory in X for action - perception associations. a)
 868 The values $x(t)$ and $g(t)$ are recorded at each t into a linked list of ordered
 869 pairs $(g(t), x(t))$, in which the links associate values $x(t - 1)$ to $g(t)$ (diagonal
 870 arrows) and $g(t)$ to $x(t)$ (within rows). Each horizontal ordered pair is an
 871 instance of the action of the composed kernel PA , during which t is constant.

872 Each diagonally-linked pair is an instance of the action of D , concurrent with
873 which t increments. b) The linked list in a) can also be represented as two
874 simple lists of ordered pairs, one representing instances of actions of D and
875 the other representing instances of actions of PA . c) The instance data in
876 either list from b) can also be represented as a matrix in which each element
877 counts the number of occurrences of an (x, g) pair. Here we illustrate just four
878 possible values of x and four possible values of g . The pair (x_1, g_1) has occurred
879 once, the pair (x_2, g_2) has occurred four times, etc. d) An RCA network that
880 constructs memories X_{MD} and X_{MPA} that count instances of actions of D and
881 PA respectively. Here X_P is the space of possible percepts and its state x_P is
882 the current percept. The space X_R is a short-term memory; its state x_R is the
883 immediately-preceding percept. The simplified notation introduced in Fig. 7 is
884 used to represent the “feedback” kernels $Copy$, M_D and M_{PA} as internal to the
885 composite RCA. The decision kernel D acts on the entire space X . The M_D
886 and M_{PA} kernels are defined in the text.

887 For the contents of memory to influence action, they must be accessible to D . They must,
888 therefore, be encoded within X . Meeting this requirement within the constraints of the CA
889 formalism requires regarding X as comprising three components, $X = X_P \times X_R \times X_M$, where
890 X_P contains percepts, X_R contains a copy of the most recent percept, and X_M contains long-
891 term memories of percept-action and action-percept associations. In this case, P becomes
892 a Markovian kernel from $W \times X_P \rightarrow X_P$ and a punctual, forgetful Markovian kernel $Copy$
893 is defined to map $X_P \rightarrow X_R$ as discussed above. The short-term memory X_R allows the
894 cross-row pairs in Fig. 8a, here written as $(x_P(t-1), g(t))$ to emphasize that $x_P(t-1)$ is a
895 percept generated by P , to each be represented as a pair $(x_R(t), g(t))$ at a single time t . To
896 be accessible to D , both these cross-row pairs and the within-row pairs $(x_P(t), g(t))$, together
897 with their occurrence counts as accumulated over multiple observations (Fig. 8c), must be
898 represented completely within X . Constructing these representations requires copying the
899 $g(t)$ components of these pairs from G to X at each t , associating the copies with either
900 $x_R(t)$ or $x_P(t)$ respectively, and accumulating the occurrence counts of the associated pairs
901 as a function of t . We define components X_{MD} and X_{MPA} of the long-term memory X_M to
902 store triples $(x_R, g_C, n_D(x_R, g_C, T))$ and $(x_P, g_C, n_{PA}(x_P, g_C, T))$ respectively, where $g_C(t)$ is
903 a copy of $g(t)$ and $n_D(x_R, g_C, T)$ and $n_{PA}(x_P, g_C, T)$ are the accumulated occurrence counts
904 of (x_R, g_C) and (x_P, g_C) , respectively, as of the accumulation time T . This T is the sum of
905 the counts stored in X_{MD} and X_{MPA} , which must be identical; the memory components
906 X_{MD} and X_{MPA} capture, in other words, the data structure of Fig. 8c completely within
907 X . To construct these memory components, we define punctual Markovian kernels $M_D : G \times X_R \times X_{MD} \rightarrow X_{MD}$ and $M_{PA} : G \times X_P \times X_{MPA} \rightarrow X_{MPA}$ (Fig. 8d) that, at each
908 t , increment $n_D(x_R, g_C, T)$ by one if x_R and g co-occur at t and increment $n_{PA}(x_P, g_C, T)$
909 by one if x_P and g co-occur at t , respectively. A similar procedure for updating “internal”
910 states on each cycle of interaction with a Markov blanket is employed in Friston (2013).
911 While we represent these memory-updating kernels as “feedback” operations in Fig. 8d
912 and in figures to follow, they can equivalently be represented as acting from G to $W \times X$

914 as in the middle part of Fig. 7.

915 The ratios $n_D(x_R, g_C, T)/T$ and $n_{PA}(x_P, g_C, T)/T$ are naturally interpreted as the frequen-
916 cies with which the pairs (x, g) have occurred as either percept-action or action-percept
917 associations, respectively, during the time of observation, i.e. between $t = 0$ and $t = T$. As
918 these values appear as components of X , they can be considered to generate, through the
919 action of some further operation depending only on X , “subjective” probabilities at $t = T$ of
920 percept-action or action-percept associations, respectively. We will abuse notation and con-
921 sider the memories X_{MD} and X_{MPA} to contain not just the occurrence counts $n_D(x_R, g_C, T)$
922 and $n_{PA}(x_P, g_C, T)$ but also the derived subjective probability distributions $\text{Prob}_D(x, g)|_{t=T}$
923 and $\text{Prob}_{PA}(x, g)|_{t=T}$ respectively. We note that these distributions $\text{Prob}_D(x, g)|_{t=T}$ and
924 $\text{Prob}_{PA}(x, g)|_{t=T}$ are subjective probabilities for the RCA encoding them, from its own in-
925 trinsic perspective. We have assumed that the kernels M_D and M_{PA} are punctual; to the
926 extent that they are not, these subjective probability distributions are likely to be inaccur-
927 ate as representations of the agent’s actual past actions and perceptions, respectively.

928 It is important to emphasize that the memory data structure shown in Fig. 8c does not
929 represent the value of the time counter t explicitly. A CA implementing this memory does
930 not, therefore, directly experience the passage of time; such a CA only experiences the cur-
931 rent values of accumulated frequencies of (x, g) pairs. However, because the current value T
932 of t appears as the denominator in calculating the subjective probabilities $\text{Prob}_D(x, g)|_{t=T}$
933 and $\text{Prob}_{PA}(x, g)|_{t=T}$, the extent to which these distributions approximate smoothness pro-
934 vides an implicit, approximate representation of elapsed time. As we discuss in §4.4 below,
935 this approximate representation of elapsed time has a natural interpretation in terms of the
936 “precision” of the memories M_D and M_{PA} , as this term is employed by Friston (2010, 2013).
937 The construction of a data structure explicitly representing goal-directed action sequences,
938 and hence the relative temporal ordering of events within such sequences, within the CA
939 framework is discussed in §4.5 below. Such a data structure is a minimal requirement for
940 directly experienced duration in the CA framework.

941 4.3 Predictive coding, goals and active inference

942 Merely writing memories is, clearly, not enough: if memories are to be useful, it must also
943 be possible to read them. Remembering previous percepts is, moreover, only useful if it
944 is possible to compare them to the current percept. As noted earlier, *exact* replication
945 of a previous percept is unlikely; hence utility in most circumstances requires *quantitative*
946 comparisons, even if these are low-resolution or approximate. These can be accomplished
947 by, for example, imposing a metric structure on X_P and all memory components computed
948 from X_P . This allows asking not just how much but in what way a current percept differs
949 from a remembered one. For now, we do this by assuming a vector space structure with
950 a norm $\|\cdot\|$ (and therefore a metric $\delta(x, x') = \|x - x'\|$) on X_P . It is also convenient to
951 assume a metric vector-space structure on G so that “similarity” between actions can be
952 discussed.

953 A vector-space structure on X_P enables talking about *components* of experience, which
 954 are naturally interpreted as basis vectors. Given a complete basis $\{\xi_i\}$ for X_P , which for
 955 simplicity is taken to be orthonormal, any percept x_P can be written as $\sum_i \alpha_i \xi_i$, where the
 956 coefficients α_i are limited to some finite resolution, and hence the vectors are limited to
 957 approximate normalization, to preserve a finite representation. The distance between two
 958 percepts $x_P = \sum_i \alpha_i \xi_i$ and $y_P = \sum_i \beta_i \xi_i$ can be defined as the distance $\delta(x_P, y_P)$.

959 To construct this vector space structure, it is useful to think of experiences in terms of
 960 “degrees of freedom” in the physicist’s sense (“macroscopic variables” or “order paramete-
 961 ters” in other literatures), i.e. in terms of properties of experience that can change in some
 962 detectable way along some one or more particular dimensions. A stationary point of light
 963 in the visual field, for example, may have degrees of freedom including apparent position,
 964 color and brightness. Describing a particular experienced state requires specifying a par-
 965 ticular value for each of these degrees of freedom; in the case of a stationary point of light,
 966 these may include x , y and z values in some spatial coordinate system and intensities I_{red} ,
 967 I_{green} and I_{blue} in a red-green-blue color space. Describing a sample of experiences requires
 968 specifying the probabilities of each value of each degree of freedom within the sample, e.g.
 969 the probabilities for each possible value of x , y , z , I_{red} , I_{green} and I_{blue} in a sample of
 970 stationary point-of-light experiences. A vector in the space X_P is then a particular combi-
 971 nation of values of the degrees of freedom that characterize the experiences in X . A basis
 972 vector ξ_i of X_P corresponds, therefore, to a particular value of one degree of freedom, e.g.
 973 a particular value $x = 1$ m or $I_{red} = 0.1$ lux. The coefficient α_i of a basis vector ξ_i is
 974 naturally interpreted as the “amount” or “extent” to which ξ_i is present in the percept;
 975 again borrowing terminology from physics, we refer to these coefficients as *amplitudes*. If
 976 α_i is the amplitude of the basis vector ξ_i representing a length of 1 m, for example, then the
 977 value of α_i represents the extent to which a percept indicates an object having a length of 1
 978 m. It is, moreover, natural to restrict the values of the amplitudes to $[0, 1]$ and to interpret
 979 the amplitude α_i of the basis vector ξ_i in the vector representation of a percept x_P as the
 980 probability that the component ξ_i contributes to x_P . This interpretation of basis vectors
 981 as representing values of degrees of freedom and amplitudes as representing probabilities is
 982 the usual interpretation for real Hilbert spaces in physics (the probability is the amplitude
 983 squared in the more typical complex Hilbert spaces).

984 The basis chosen for X_P determines the bases for X_R , X_{MD} and X_{MPA} . It must, moreover,
 985 be assumed that elements of these latter components of X are experientially tagged as such.
 986 An element x_R in X_R must, for example, be experienced differently from the element x_P in
 987 X_P of which it is a copy; without such an experiential difference, previous, i.e. remembered
 988 and current percepts cannot be distinguished as such from the intrinsic perspective. The
 989 existence of such experiential “tags” distinguishing memory components is a prediction of
 990 the current approach, which places all memory components on which decisions implemented
 991 by D can depend within the space X of experiences. Models in which some or all compo-
 992 nents of memory are implicit, e.g. encoded in the structure of a decision operator, require
 993 no such experiential tags for the implicit components. It is interesting in this regard that
 994 humans experientially distinguish between perception and imagination (a memory-driven

995 function), that this “reality monitoring” capability appears to be highly but not exclusively
 996 localized to rostral prefrontal cortex, and that disruption of this capability correlates with
 997 psychosis (Simons, Gilbert, Henson and Fletcher, 2008; Burgess and Wu, 2013; Cannon,
 998 2015). Humans also experientially distinguish short-term “working” memories from long-
 999 term memories. We predict that specific monitoring capabilities provide the experiential
 1000 distinctions between short- (e.g. X_R) and long-term (e.g. X_{MD} and X_{MPA}) memories and
 1001 distinguish functionally-distinct long-term memory components from each other. From a
 1002 formal standpoint, such distinguishing tags can be considered to be additional elements in
 1003 each vector in each of the derived vector spaces; while such tags play no explicit role in the
 1004 processing described below, their existence will be assumed.

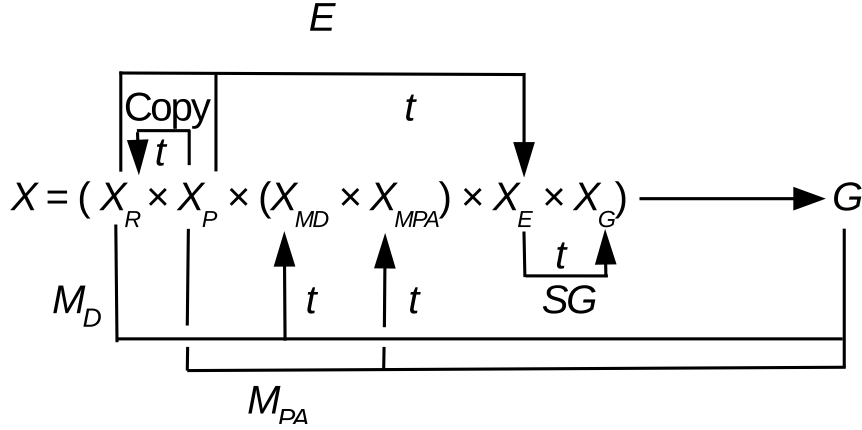
1005 As the memories X_{MD} and X_{MPA} and hence the conditional probability distributions
 1006 $\text{Prob}_D(x(t), g(t)|x(t-1), g(t-1))$ and $\text{Prob}_{PA}(x(t), g(t)|x(t-1), g(t-1))$ contain informa-
 1007 tion about the observer’s entire experience of the world, they enable differential responses
 1008 to $x_R - g$ or $g - x_P$ pairings that evoke different degrees of “surprise” by either confirming
 1009 or disconfirming previous associations to different extents. We note that the term ‘surprise’
 1010 is being used here in its informal sense of an *experienced* departure from expectations, not
 1011 in the technical sense employed by Friston (2010; 2013; see also Friston et al., 2015; Fris-
 1012 ton et al., 2016) to refer to an event that causes or threatens to cause a departure from
 1013 homeostasis and hence has negative consequences for fitness. To implement such differen-
 1014 tial responses to surprise, it is natural to choose functions for updating these conditional
 1015 probability distributions that depend on the vector distance(s) between the percept x_R (for
 1016 $\text{Prob}_D(x(t), g(t)|x(t-1), g(t-1))$) or x_P (for $\text{Prob}_{PA}(x(t), g(t)|x(t-1), g(t-1))$) and the
 1017 percept(s) previously associated, within X_{MD} and X_{MPA} respectively, with g . Functions
 1018 can clearly chosen that either enhance or suppress memories of surprising events. This
 1019 generalization requires no additional components or elements within X ; hence it enhances
 1020 function without altering the architecture.

1021 The simplest possible action is no action: the agent merely observes the world. The extremal
 1022 outcomes of such observation are on the one hand James’ “blooming, buzzing confusion,” i.e.
 1023 a completely random $x_P(t)$, and on the other stasis, a fixed and invariant $x_P(t)$. Memory is
 1024 obviously useless in either case; indeed, the latter corresponds to the “dark room” situation
 1025 discussed above. Memory becomes useful if a world on which no action is taken generates
 1026 some number of the possible percepts significantly more often than the others. The same
 1027 is true in the case of any other constantly-repeated action. It is equivalent to say: any
 1028 action which, when repeated indefinitely, is followed by either random or static percepts
 1029 is a useless action to take. Such an action has no “epistemic value” in the sense used by
 1030 Friston et al. (2015). Randomness and stasis may be useful as *components* of experience -
 1031 indeed as discussed below, stasis is a *necessary* component of useful experience - but only
 1032 when embedded in non-random, non-static contexts. Let us assume, therefore, that RCAs
 1033 of interest are embedded in W s that generate non-random, non-static percepts in response
 1034 to all actions. Note that this assumption is consistent with ITP: it does not require either
 1035 P or A to respect the causal structure of W .

1036 In a non-random, non-static world, the memories X_{MD} and X_{MPA} provide a basis for

1037 predictive coding: the probability assigned to an action g at $t + 1$ can depend on the vector
 1038 difference between the current percept $x_P(t)$ and previous percepts either immediately-
 1039 antecedent or immediately-consequence to actions like g . A percept $x_P(t)$ can, in this case,
 1040 “predict” an action $g(t + 1)$ that is “expected,” on the basis of the probabilities stored
 1041 in X_{MPA} , to result in a subsequent percept $x_P(t + 1)$ that is either similar or dissimilar
 1042 to $x_P(t)$. Assigning high probabilities to actions at $t + 1$ expected to result in percepts
 1043 similar to $x_P(t)$ is implicitly “evaluating” $x_P(t)$ as in some sense “good” or “desirable,”
 1044 while assigning low probabilities to actions at $t + 1$ expected to result in percepts similar to
 1045 $x_P(t)$ is implicitly evaluating $x_P(t)$ as in some sense bad or undesirable. These operational
 1046 senses of “good” and “bad” percepts are consistent with the senses of “good” and “bad”
 1047 percepts as enhancing or threatening the maintenance of homeostasis employed by Friston
 1048 (2010; 2013). A “bad” experience in this operational sense is a outcome that an agent
 1049 did not expect to experience, i.e. a stressor such as being hungry or poor, on the basis
 1050 of the implicit “model” of W encoded by the probability distributions contained in the
 1051 memories X_{MD} and X_{MPA} . In the limit, a maximally “bad” experience is one that violates
 1052 the fundamental expectation that experiences will continue that is encoded by all non-
 1053 zero values of the subjective probabilities $\text{Prob}_D(x, g)|_{t=T}$ and $\text{Prob}_{PA}(x, g)|_{t=T}$; such an
 1054 experience destroys connectivity between the agent in question and the surrounding RCA
 1055 network (i.e. the agent’s W), setting the agent’s fitness to zero and corresponding to the
 1056 “death” of the agent as discussed in §3.3 above.

1057 This evaluative function can be made explicit by representing it as a distinct operation. To
 1058 do this, we add a further memory component X_E to X . To allow for the possibility that
 1059 an observer has “innate” biases toward or against particular percepts, we consider X_E to
 1060 comprise two probability distributions, $\text{Prob}_{good}(x_P)$ and $\text{Prob}_{bad}(x_P)$, with *a priori* values
 1061 fixed at $t = 0$. Such innate evaluation biases can be considered to be innate “preferences”
 1062 or “beliefs” as they often are in the infant-cognition literature (e.g. Baillargeon, 2008;
 1063 Watson, Robbins and Best, 2014). We represent the evaluation operation E as having two
 1064 components $E = (E_{good}, E_{bad})$, where E_{good} is a punctual kernel $X_P \times X_R \times E \rightarrow E$ that
 1065 updates $\text{Prob}_{good}(x_P)$ at each t and E_{bad} is a punctual kernel $X_P \times X_R \times X_E \rightarrow X_E$ that
 1066 updates $\text{Prob}_{bad}(x_P)$ at each t . For simplicity, we assume that E_{good} increases $\text{Prob}_{good}(x_P)$
 1067 by a factor ≥ 1 that approaches unity as $\text{Prob}_{good}(x_P) \rightarrow 1$ whenever both $\text{Prob}_{good}(x_P(t)) >$
 1068 0 and $\text{Prob}_{good}(x_R(t)) > 0$ and that E_{bad} increases $\text{Prob}_{bad}(x_P)$ by a factor with similar
 1069 behavior whenever both $\text{Prob}_{bad}(x_P(t)) > 0$ and $\text{Prob}_{bad}(x_R(t)) > 0$. This E effectively
 1070 implements the heuristic: an experience is remembered as better if it is followed by a good
 1071 experience, and remembered as worse if it is followed by a bad experience. Note that while
 1072 this heuristic is consistent with the association of “good” and “bad” with maintaining or
 1073 not maintaining either homeostasis or connectivity as discussed above, it also allows a
 1074 given x_P to be both probably good and probably bad, a not-unrealistic situation. This
 1075 additional structure on X is summarized in Fig. 9. Extending the evaluative process from
 1076 the scalar representation provided by these probabilities to a multidimensional, i.e. vector,
 1077 representation costs memory and kernel complexity but does not change the architecture.



1078 *Fig. 9:* Adding memories for evaluations of percepts (X_E) and for a current
 1079 goal (X_G) to Fig. 7d. Connections to W have been elided for clarity.

1080 Evaluating percepts implicitly evaluates the actions that are followed by those percepts;
 1081 this implicit transfer of estimated “good” or “bad” value from percepts to actions is now
 1082 implemented by D . A “rational” D , for example, would assign high probabilities to actions
 1083 g that are associated in X_{MPA} with subsequent percepts that have high valuations in X_E .
 1084 If W is such that the relative ranking of percepts by value changes only slowly with t ,
 1085 relatively highly- and lowly-ranked percepts can be considered to be positive and negative
 1086 “goals” respectively. As Friston (2010, 2013) has emphasized, goals are effectively long-term
 1087 expectations to which an uncertainty-minimizing agent attempts to match perceptions;
 1088 Friston and colleagues call acting so as to match perceptions to goals “active inference.”
 1089 Within the CA framework, the minimal functional architecture required for active inference
 1090 is that shown in Fig. 9. Here a memory component X_G holds the current goal; it is
 1091 populated by a punctual, forgetful kernel SG acting on X_E . While SG can be taken to
 1092 choose percepts of high value as goals, its specific action can be left open. Note than in this
 1093 architecture, incremental adjustments of the “world model” X_{MPA} and “self model” X_D
 1094 are made in parallel with active inference: expectations are modified to fit perceptions even
 1095 when actions are taken to modify perceptions to fit expectations. Note also that placing
 1096 the evaluation and goal memories X_E and X_G within the experience space X is predicting
 1097 that the contents of these memories are both experienced and experienced as distinct, as
 1098 they indeed are in neurotypical humans. While the specific mechanisms implementing the
 1099 experiential distinction between these memory components remains uncharacterized, the
 1100 present framework predicts that such mechanisms exist.

1101 By iteratively constructing representations of the antecedents and consequences of actions,
 1102 the kernels M_D and M_{PA} implement a simple kind of learning. The operator E similarly

1103 implements a simple form of evaluative feedback. The action choices made by D can,
1104 therefore, progressively improve with experience. It is important to emphasize that M_D ,
1105 M_{PA} , E , SG and D are all by assumption homogeneous kernels. What changes as the
1106 system learns is not the choice function D , but the contents of the data structures – the
1107 memories X_{MD} , X_{MPA} , X_E and X_G – that serve as ancillary inputs to D . The “knowledge”
1108 of an RCA with this architecture is, therefore, entirely explicit. This is marked contrast
1109 to typical neural-network models, including recent “deep learning” models (for a recent
1110 review, see Schmidhuber, 2015), in which learning is entirely implicit and the decision rules
1111 learned are notoriously hard to reverse engineer. It is worth noting that standard neural-
1112 network models have no intrinsic perspective; as emphasized earlier, it is the requirement
1113 that an RCA learns about W from its own intrinsic perspective that forces what is learned
1114 to be made explicit in a memory located in X , i.e. in a memory encoding contents that
1115 are experienced - but are not necessarily reportable - by the RCA. While the kernels M_D ,
1116 M_{PA} , E , SG , as well as others to be introduced below, that populate explicit memories
1117 can, together with the decision kernel D be considered to encode implicit memories in the
1118 current model, the assumption that all such kernels are homogeneous implies that these
1119 implicit memories are not loci of learning. The kinds of “practised skill” memories that
1120 are canonically regarded as implicit are most naturally modelled as structures, e.g. fixed
1121 or fully-automatized learned action patterns, within the action space G in the current
1122 framework; an exploration of such structures are developed within G is beyond the present
1123 scope.

1124 It is important to note that whether D is “rational” in the sense of favoring actions that re-
1125 sult in “good” outcomes, and hence the extent to which the choices favored by D “improve”
1126 with experience, is left open within the architecture. If W is such that “good” choices cor-
1127 relate with the acquisition of resources required for survival, a basic orientation or “drive”
1128 toward increasing the average subjective valuation of “good” percepts can be expected to
1129 emerge in a population of agents whenever the required resources are scarce. Friston has
1130 argued that predictability of experience is itself the primary resource that organisms seek
1131 to maximize, and that the drive to pursue and acquire external resources can be under-
1132 stood in terms of maintaining the predictability of experiences that facilitate or enhance
1133 the maintenance of physiological homeostasis (Friston, 2010; 2013; Friston, Thornton and
1134 Clark, 2012). Reducing the uncertainty of experiences from a large environment requires
1135 extensive sampling of the environment’s behaviors and hence active exploration; effective
1136 agents in a large W can, therefore, be expected to display a “curious rationality” that
1137 maintains homeostasis while devoting significant energy to active exploration and learning
1138 (reviewed by Gottlieb, Oudeyer, Lopes and Baranes, 2013). Friston et al. (2015; 2016)
1139 make a similar point: the minimization of expected surprise in the strict sense of departure
1140 from homeostasis (i.e. the minimization of variational free energy) contingent upon remem-
1141 bered action-perception associations can always be expressed as a mixture of “epistemic”
1142 and “pragmatic” value. The pragmatic value is the expected outcome according to prior
1143 preferences, i.e. “good” or “bad” evaluations, while the epistemic value is the utility of the
1144 action for learning, i.e. reducing the potential for uncertainty or surprise in the future. This
1145 resolution of uncertainty through active sampling is at the heart of many active inference

1146 schemes and arises naturally in any model in which the agent expects to occupy the states
1147 it prefers.

1148 4.4 Reference frames and attention

1149 While defining expectations over percepts can be expected to be useful in some circum-
1150 stances, many aspects of realistic behavior require defining and acting on expectations
1151 defined over individual or small subsets of *components* of percepts. The memories X_{MD}
1152 and X_{MPA} together provide the data needed to allow individual component - action as-
1153 sociations to be computed; the memory X_E similarly provides the data needed to allow
1154 individual component valuations to be computed. Let X_C and X_{EC} be memories that store
1155 conditional probability distributions and evaluations, respectively, of individual components
1156 of percepts. To define X_C , note that the $x_R - g$ and $g - x_P$ associations stored in X_{MD} and
1157 X_{MPA} respectively allow each action g to be viewed as a relation $\{(x_R, x_P)\}$ implemented
1158 by PA . Expressing these percepts as vectors $x_R(t) = \sum_i \alpha_i(t)\xi_i$ and $x_P(t) = \sum_i \beta_i(t)\xi_i$,
1159 we can view the action of g on the component ξ_i at t as $g_{\xi_i}(t) : \alpha_i(t) \mapsto \beta_i(t)$. Each g
1160 can, in other words, be viewed as increasing or decreasing the amplitude of each percep-
1161 tual component ξ_i from one percept to the next. As it is natural to view amplitudes as
1162 probabilities of occurrence as discussed above, each g can be viewed as increasing or de-
1163 creasing the probability of each perceptual component ξ_i from one percept (i.e. value of
1164 t) to the next. The memory X_C can, therefore, be viewed as storing t -indexed conditional
1165 probabilities $\text{Prob}_t(\xi_i|g, \text{Prob}_{t-1}(\xi_i))$ of perceptual components given actions. To update
1166 the distribution of $\text{Prob}_t(\xi_i|g, \text{Prob}_{t-1}(\xi_i))$ as a function of t , we define a punctual kernel
1167 C as a map $X_{MD} \times X_{MPA} \times X_C \rightarrow X_C$. Subject to the constraint that all probabilities
1168 remain normalized, this map can in principle implement any arbitrary updating function.

1169 The memory X_{EC} containing component valuations may be constructed from X_E in a sim-
1170 ilar fashion, by defining punctual, forgetful kernels EC_{good} and EC_{bad} that map $X_E \rightarrow$
1171 X_{EC} . The kernels EC_{good} and EC_{bad} assign, respectively, “good” valuations to components
1172 strongly represented in “good” percepts and “bad” valuations to components strongly rep-
1173 resented in “bad” percepts. A suitable function for each would assign to each component
1174 ξ_i the average valuation of percepts x_P in which the coefficient α_i of ξ_i is greater than
1175 some specified threshold. With additional memory, this mechanism can be extended to
1176 assign values to (finite ranges of) amplitude values of components. Note that component
1177 valuations constructed in this way are in an important sense context-free; representing com-
1178 ponent valuations conditioned on the valuations of other components requires both more
1179 memory and more complex kernels.

1180 The memory components X_C and X_{EC} provide the “background knowledge” required for
1181 component-directed as opposed to entire-percept directed actions. What remains to be
1182 constructed is a process of selecting a component on which to act, and a second component
1183 with respect to which the action is taken. Consonant with current usage in physics (e.g.
1184 Bartlett, Rudolph and Spekkens, 2007), we refer to this second, context-setting component
1185 as a reference frame for the action. Specifying a reference frame is specifying what does

1186 *not* change when an action is taken; hence reference frames provide the basis for specifying
1187 what does change. Reference frames provide, in other words, the necessary stasis with
1188 respect to which change is perceptible. Measurement devices such as meter sticks provide
1189 the canonical example: a measurement made with a meter stick is only meaningful if one
1190 assumes that the actions involved in making the measurement do not change the length
1191 of a meter stick. More broadly, any context in which observations are made, whether a
1192 particular laboratory set-up or an everyday scene, is meaningful as a context only if it
1193 itself does change as a result of making the observation. A reference frame is, therefore, a
1194 *stipulated* solution to the frame problem, the problem of specifying what does not change
1195 as a result of an action (McCarthy and Hayes, 1969; reviewed by Fields, 2013b). Such
1196 stipulations are inherently fragile and defeasible: a context that does observably change,
1197 like a “meter stick” with an observably context-dependent length, ceases to be a reference
1198 frame as soon as its variation is detected. Stipulated reference frames are, nonetheless,
1199 *useful* solutions to the frame problem to the extent that they enable successful behavior in
1200 the niche of the agent employing them. Absent a level of control over the environment that
1201 ITP forbids, they are the only kinds of reference frames available.

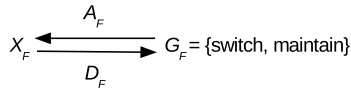
1202 While the frame problem has a long history in AI, its impact on cognitive science more
1203 generally has been primarily philosophical (see, e.g. the contributions to Pylyshyn (1987)
1204 and Ford and Pylyshyn (1996)). The question of how human perceivers identify *contexts*
1205 as opposed to objects or events and how they detect changes in context have received little
1206 direct investigation. The current model predicts that contexts are defined constructively
1207 by the activation of discrete reference frames that impose expectations of constancy and
1208 limit attention to features expected to remain constant. Experimental demonstrations of
1209 change-blindness (reviewed by Simons and Ambinder, 2005) show that such limitations of
1210 attention exist. Virtual reality methods provide opportunities to experimentally manipulate
1211 context identification, and hence to probe the specific reference frames employed to identify
1212 contexts, in ways that remain largely unexplored.

1213 For complex organisms, the most important reference frame is arguably the *experienced self*,
1214 generally including one or more distinguishable components of the *body*. This experienced
1215 self reference frame comprises a collection of components of experience that do not change
1216 during some, most or even all actions. The experienced self as a reference frame appears to
1217 be innate in humans (e.g. Rochat, 2012) and may be innate in higher animals generally. It is
1218 with respect to the experienced self as a reference frame that infants learn their capabilities
1219 for actions as bodily motions and for social interactions as communications with others (e.g.
1220 von Hofsten, 2007). Actions of or on the body, e.g. moving a limb, require that other parts
1221 of the experienced self, e.g. the mass and shape of the limb and its point of connection
1222 to the rest of the body, remain fixed to serve as the reference frame for the action. As
1223 the body grows and develops, its representation must be updated to compensate for these
1224 changes if its function as a reference frame is to be preserved. The experienced self reference
1225 frame is readily extensible to tools, vehicles, and fully-virtual avatars in telepresence and
1226 virtual-reality applications, and is readily manipulated in the laboratory. Disruptions of the
1227 experienced self as a reference frame present as pathologies ranging from schizophrenia to

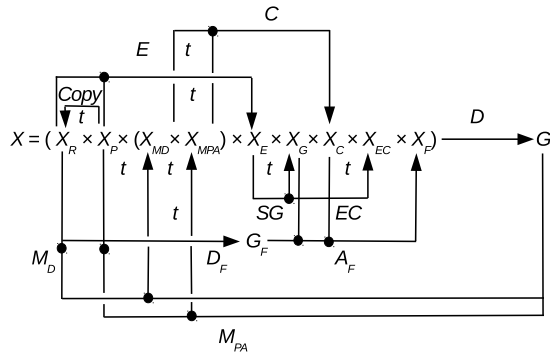
1228 anosognosia. These latter provide a clinical window into the human implementation of the
 1229 bodily and emotive self as a fusion of interoceptive and perceptual inputs (e.g. Craig, 2010;
 1230 Seth, 2013) and of the cognitive self as a fusion of memory-access and executive functions
 1231 that develops gradually from infancy to early adulthood (e.g. Simons, Henson, Gilbert and
 1232 Fletcher, 2008; Metzinger, 2011; Hohwy, 2016).

1233 Selecting a particular component of a percept on which to act and another component
 1234 or components, such as the experienced self or the experienced self in some perceived
 1235 surroundings, to serve as a fixed context for an action is an act of *attention*. The selected
 1236 components must, moreover, remain subjects of attention throughout the action. Any
 1237 agent capable of attending to some component of an ongoing scene must also, however,
 1238 be capable of switching attention to a different component if something unexpected and
 1239 important happens. Attention requires, therefore, not just a decision about what to attend
 1240 to, but also a decision about whether to maintain or switch attentional focus. To meet these
 1241 requirements, we introduce an “attentional workspace” X_F , a memory that contains a goal-
 1242 dependent focus of attention ξ_i , a focus-dependent reference frame ξ_j and a time counter
 1243 t_F that measures the duration of an attentional episode. We also define an attentional
 1244 action space G_F containing two actions, ‘switch’ and ‘maintain’ that alter or preserve the
 1245 attentional focus, respectively, and a forgetful punctual kernel $D_F : X_P \times X_R \times X_E \times X_G \rightarrow$
 1246 G_F that selects $g_F =$ ‘switch’ at t if the valuation of $x_P(t)$ differs from that of $x_R(t)$ by
 1247 some specified threshold and selects $g_F =$ ‘maintain’ otherwise. These elements of G_F
 1248 correspond to actions A_F on the workspace X_F , as shown in Fig. 10a. The action A_{Fm}
 1249 selected by $g_F =$ ‘maintain’ only increments t_F . The action A_{Fs} selected by $g_F =$ ‘switch’
 1250 selects a new focus of attention ξ_k , a new reference frame ξ_l and resets t_F to zero. We
 1251 represent this action as a forgetful punctual kernel $A_{Fs} : X_P \times X_G \times X_C \times X_{CE} \rightarrow X_F$.
 1252 How this attention-switching kernel is defined has a potentially large impact on the behavior
 1253 of the RCA whose attentional workspace X_F it affects. A rational A_{Fs} could be expected
 1254 to select a component ξ_i on which to focus that had a relatively large amplitude α_i in
 1255 both the current percept x_P and a high-value goal and a reference frame ξ_j , also with a
 1256 relatively large amplitude in both x_P and the goal, that was affected in the past primarily
 1257 by actions that did not affect ξ_i . While the valuation of the attentional focus ξ_i may be
 1258 “bad,” a rational A_{Fs} would select a reference frame ξ_j with a “good” or at least not “bad”
 1259 valuation, as this amplitude of this component is meant to be kept fixed in subsequent
 1260 interactions with W . A rational D kernel acting on the workspace X_F would then choose
 1261 actions g that, in the past as recorded in X_C , moved the amplitude of x_i in the direction of
 1262 its value in the chosen goal state while keeping the amplitude of x_j fixed. As X_C , X_{EC} and
 1263 X_F are updated one cycle behind X_{MD} , X_{MPA} , X_E and X_G and hence two cycles behind
 1264 X_P , the kernel D must always work with expectation and valuation information that is
 1265 slightly out-of-date.

a)



b)



1266 *Fig. 10:* a) Kernels that maintain or switch attentional focus. b) Additions to
 1267 *Fig. 9* required to support attention. Connections to W are again elided for
 1268 clarity.

1269 The structure of and operations within the experiential space X required for an atten-
 1270 tional system are summarized in *Fig. 10b*. Selecting a new component for attention and
 1271 maintaining attention on a previously-selected component are competitive processes in this
 1272 architecture, as they are in humans (reviewed by Vossel, Geng and Fink, 2014). When
 1273 top-down goals and expectations dominate and hence the dorsal attention system controls
 1274 perceptual processing, the salience of goal-irrelevant stimuli is reduced; a switch to vigilance
 1275 and hence ventral attentional control, in contrast, reduces the salience of goal-relevant stim-
 1276 uli. Top-down, dorsal attentional dominance facilitates exploration and information gather-
 1277 ing, while bottom-up, ventral attentional dominance facilitates threat avoidance. This
 1278 attention switch can be incorporated into predictive coding and active inference models
 1279 using the concept of “precision” for both expectations and percepts; high-precision expect-
 1280 ations dominate low-precision percepts and vice-versa (Friston, 2010; 2013). Precision is
 1281 effectively a measure of reliability based on prior experiences and is hence a second-order
 1282 expectation that must be learned by refining an *a priori* bias as discussed above. Predic-
 1283 tive coding networks modulated by estimated precision have been shown to describe the
 1284 cellular-scale connection architecture of cortical minicolumns (Bastos et al., 2012) as well
 1285 as the modular connection architectures of motor (Shipp, Adams and Friston, 2013) and
 1286 visual (Kanai, Komura, Shipp and Friston, 2015) processing (see also Adams, Friston and
 1287 Bastos (2015) for an overview of these results). As noted earlier, the smoothness of stored
 1288 probability distributions provides a natural estimate of the number of experiences that have
 1289 contributed to them and hence their reliability. A rational switching function can be ex-
 1290 pected to favor high-reliability expectations and disfavor low-reliability expectations, and
 1291 hence to implement a precision-based modulation of attention.

1292 Extending the system shown in *Fig. 10b* to multiple focus and/or reference components
 1293 costs memory and processing complexity, but does not change the architecture. It is inter-

1294 esting to note that within this architecture, all change is implicitly attributed by the agent
1295 to the action taken; from the agent’s intrinsic perspective, its actions change the state of its
1296 attentional focus with respect to its reference frame. For the system to behave effectively,
1297 the world W must be such that this attribution of observed changes to executed actions is
1298 satisficing in W . The world must not, in other words, surprise the agent so often that the
1299 agent’s sense that actions have predictable consequences becomes impossible to maintain.
1300 The world must not, in other words, exhibit either overall randomness or overall stasis as
1301 noted earlier.

1302 It is worth re-emphasizing, moreover, that in the CA framework X is a space of *experi-*
1303 *ences*. Hence the RCA depicted in Fig. 10b is regarded as *experiencing* each state of its
1304 highly-structured space X , including all those components on which its attention is *not*
1305 focussed (the formalism leaves open the question of whether these components themselves
1306 have unexperienced internal structure). It may, however, be “unconscious” of unattended
1307 components in the sense in which this term is used in theories that associate consciousness
1308 with relative amplification or attention (e.g. Baars, Franklin and Ramsay, 2013; Dehaene,
1309 Charles, King and Marti, 2014; Graziano, 2014). In general, how an RCA acts depends
1310 on its attentional focus. Reporting what it is experiencing, e.g. to an investigator in a
1311 laboratory or even to itself via a modality such as inner speech, is a specific kind of ac-
1312 tion that requires a specific attentional focus. Whether the attentional focus required to
1313 support a given form of reporting is achieved in any particular case or is even achievable
1314 by a particular RCA is a matter of architecture, i.e. of how the memory-construction and
1315 attentional-control kernels are defined. Agents that never report particular kinds of experi-
1316 ences, or that never report experiences using a given modality such as inner speech (Heavey
1317 and Hurlburt, 2008), are not only possible but to be expected within the CA framework.
1318 Indeed the CA framework predicts that *agents are typically aware of more than they can*
1319 *report awareness of* to an external observer or even to themselves. Agents are, in other
1320 words, typically *under-equipped with attentional resources*, and hence unable to access some
1321 or even much of their experience for behavioral reporting via any particular modality. Being
1322 under-equipped for reporting experiences *post hoc* is unsurprising on evolutionary grounds;
1323 indeed why human beings should engage in so much *post hoc* self-reporting via modalities
1324 such as inner speech remains a mystery (Fields, 2002). As reportability by some observable
1325 behavior remains the “gold standard” in assessments of awareness (e.g. Dehaene, Charles,
1326 King and Marti, 2014), this strong and counter-intuitive prediction of the CA framework
1327 can at present only be tested indirectly, e.g. using phenomena such as blindsight (re-
1328 viewed by Overgaard, 2011). It raises the methodological question of whether “reporting”
1329 of experiences by imaging methods such as fMRI, as employed by Boly, Sanders, Mashour
1330 and Laureys (2013), for example, with otherwise-unresponsive coma patients, should be
1331 regarded as evidence of awareness across the board.

4.5 Remembering and planning action sequences

The attentional workspace X_F defined above does not explicitly represent the action taken at each t and so cannot support either memory for “cases” of successful action or planning. The most recently executed g is, however, available within X_{MD} . A fixed-capacity case memory can be regarded as a subjective probability distribution over possible cases, where each case is a vector of fixed length l_{case} , the components of which are quadruples $(\alpha_i \xi_i, \beta_j \xi_j, t_F, g(t_F))$ with the percept components ξ_i, ξ_j and the amplitude β_j fixed. A case defined in this way provides a representation of how the amplitude α_i of the attentional focus ξ_i varies relative to the fixed amplitude β_j of the reference frame ξ_j when subjected to the sequence $g(t_F = 0) \dots g(t_F = l_{case})$ of actions. This definition formulates in language compliant with ITP the concept of a case employed in the case-based reasoning and planning literature (Riesbeck and Schank, 1989; Kolodner, 1993). It is also similar in both role and scope to the concept of an “event file” introduced by Hommel (2004) to represent the temporal binding of perceptions with context-appropriate actions. Cases or event files are effectively “snapshots” of active inference that show how a particular perceptual input is processed given the attentional context in which it is received and the particular expectations that it activates.

As an example, consider a sequence of actions involved in reaching for and grasping a coffee cup. The immediate goal of the sequence is to grasp the coffee cup; we will ignore the question of different grasps being needed for different subsequent actions. The target of the sequence is a *particular* coffee cup that is visually identifiable by particular perceived features, e.g. location, size, shape and color. The cup’s perceived size, shape and color do not change as a result of the motion; hence their values can serve as the reference frame that determines the cup’s identity. As the goal of the action sequence is to change the perceived location of the coffee cup, its location cannot be included in the reference frame; if it was, the cup would lose its identity when it was moved. The attentional workspace X_F , therefore, contains the variable perceived values of the positions of the cup and of the reaching hand as foci and the fixed perceived values of the size, shape and color of the cup as the reference frame. The recorded case contains, effectively, a sequence of “snapshots” of the contents of X_F : a time sequence of cup and hand position values, together with the actions that produced them, relative to these fixed reference values. A memory M_{case} for such cases can be constructed using the counter-incrementing methods used to construct X_{MD} and X_{MPA} above. As action sequences that are worth recording are typically those that either satisfied goals or led to trouble, it is useful to construct each record in M_{case} as a 5-tuple $[x_P(t_F = 0), E((x_P(t_F = 0))), x_P(t_F = l_{case}), E((x_P(t_F = l_{case}))), case(t_F)]$, where $x_P(t_F = 0)$ and $x_P(t_F = l_{case})$ are the full percepts at the beginning and the end of $case(t_F)$ respectively, and $E((x_P(t_F = 0)))$ and $E((x_P(t_F = l_{case})))$ are their evaluations as recorded in X_E . This representation allows M_{case} to be searched – i.e. kernels acting on M_{case} to depend upon – either the initial state and its evaluation or the final state and its evaluation. Case memories constructed in this way are clearly combinatorially explosive; hence case-based planning in systems with limited memory is necessarily heuristic, not exhaustive, a condition widely recognized in the case-based planning literature.

1374 It is natural to interpret a set of one or more fixed components of experience, with respect
1375 to which one or more other components of experience change when one or more sequences
1376 of actions is executed as defining an effective or apparent *object*. Objects defined in this way
1377 are collections of expectations, based on accumulated experience, about the co-occurrence
1378 and co-variation under actions of particular values of particular experiential degrees of
1379 freedom. Objects in this sense are effectively *categories* defined by fixed (i.e. reference) and
1380 variable features together with sets of expected behaviors, i.e. changes in the amplitudes of
1381 the variable features relative to the fixed features in response to actions. Hence such objects
1382 are more properly considered to be object *types* as opposed to *de re* individuals. While
1383 an agent may *assume*, as a useful heuristic, that an object category has only one member
1384 and act on the basis of this assumption, consistency with ITP requires that nothing in
1385 the agent’s experience can be sufficient to demonstrate that this is the case. Hence object
1386 identity over time is ambiguous in principle in the ITP/CA framework. Objects defined in
1387 this way play the role of “icons” on the ITP interface. As the number of recorded cases
1388 involving actions that change the state of some object increase, its “icon” gains predictable
1389 functionality and hence utility as a locus of behavior.

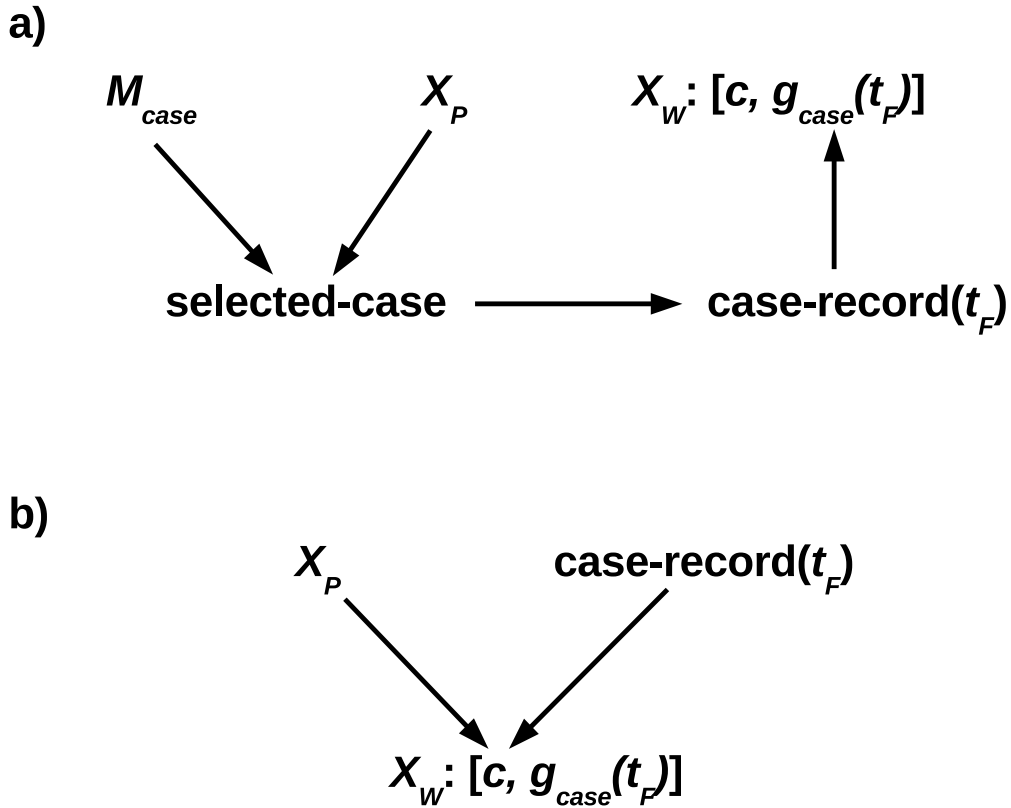
1390 The present framework leaves open the question of whether any “object”-specifying ref-
1391 erence frames are innate. It predicts, however, that any such reference frames, whether
1392 innately specified or constructed from experience, will have low dimensionality compared
1393 to the perceptual experiences that they help to interpret. Dramatic evidence for low di-
1394 mensionality is provided by studies of two of the earliest-developing and ecologically most
1395 crucial reference frames for humans, those that identify animacy and agency (reviewed by
1396 Scholl and Tremoulet, 2000; Scholl and Gao, 2013; Fields, 2014). Indeed Gao, McCarthy
1397 and Scholl (2010) have shown that a simple oriented “V” shape not only satisfies the typical
1398 human visual criterion for agency detection, but distracts attention sufficiently to disrupt
1399 performance in an object-tracking task. Human face-recognition criteria are similarly rudi-
1400 mentary. Additional evidence for low reference-frame dimensionality is provided by the
1401 kinds of categorization conflicts studied in the quantum cognition literature (reviewed e.g.
1402 by Pothos and Busemeyer, 2013; Bruza, Kitto, Ramm and Sitbon, 2015), for example the
1403 “Linda” problem. Here the “natural” reference frames, i.e. concepts or coherent sets of
1404 expectations, do not exhibit classical compositionality; combining reference frames to repro-
1405 duce the judgements made by subjects requires the use of complex “quantum” probability
1406 amplitudes. Complex probabilities can, however, be represented by classical probabilities
1407 in higher-dimensional spaces (e.g. Fuchs and Schack, 2013; see also Fields, 2016 for a less
1408 formal discussion), consistent with attentional selection of a low-dimensional subspace to
1409 serve as a reference frame. If “object”-specifying reference frames in fact encode fitness
1410 information as ITP requires, one would expect a general inverse correlation between fitness
1411 consequences and reference frame dimensionality. While both the global and local struc-
1412 ture of the typical human category hierarchy have been investigated (reviewed by Martin,
1413 2007; Keifer and Pulvermüller, 2012), neither the minimal functional content (i.e. dimen-
1414 sionality) nor the fitness-dimensionality correlation of typical categories have been broadly
1415 investigated.

1416 The components of the experienced self reference frame, taken together, constitute an
 1417 iconic object – the experienced self as a persistent embodied actor – in the above sense.
 1418 The features of the experienced self as persistent embodied actor that are employed as
 1419 fixed reference features with respect to which other features of the experienced self are
 1420 allowed to vary change only slowly and asynchronously as a function of time; it is this slow
 1421 and asynchronous change in reference features that allow the approximation of a persistent
 1422 experienced self (but see Klein, 2014 for a discussion of the sense of a persistent experienced
 1423 self in the presence of conflicting perceptual evidence). The conditions under which non-
 1424 self objects are represented as persistent over extended time, in particular across extended
 1425 periods of non-observation, have been subjected to surprisingly little direct experimental
 1426 investigation and are not well understood (e.g. Scholl, 2007; Fields, 2012). Both the
 1427 extensibility of the experienced self reference frame to incorporate otherwise non-self objects
 1428 discussed earlier and the sheer variety of pathologies of the experienced self, including
 1429 depersonalization syndromes (e.g. Debruyne, Portzky, Van den Eynde and Audenaert,
 1430 2009), suggest that the experienced self - non-self distinction is not constant for individual
 1431 human subjects and highly variable between subjects. This question cannot, unfortunately,
 1432 yet be addressed productively in non-human subjects.

1433 With this concept of an iconic object, the functional difference between a case memory
 1434 M_{case} and the event memories X_{MD} and X_{MPA} becomes clear: M_{case} records sequences of
 1435 *partial* events in which, in each sequence, only the response to actions of the attentional
 1436 focus ξ_i and the lack of response to actions of the reference ξ_j are made explicit. Each case
 1437 in M_{case} can, therefore, be thought of as imposing an implicit, goal-dependent criterion of
 1438 *relevance* on the actions it records.

1439 Recording object-directed action sequences is useful to an agent because it enables previously-
 1440 successful sequences to be repeated and previously-unsuccessful sequences to be avoided.
 1441 Selecting a previously-recorded case from memory for execution under some similar cir-
 1442 cumstances is the simplest form of planning. Executing the action sequence recorded in a
 1443 remembered case requires, however, shortcutting the usual decision process D . Within the
 1444 architecture shown in Fig. 10, the simplest way to accomplish this is to associate a working
 1445 memory X_W with the attentional focus X_F , and to include in X_W a control bit c on which
 1446 D depends. If $c = 0$, D is independent of the contents of X_W and acts as in Fig. 9. If $c = 1$,
 1447 D selects the action g represented in X_W . Populating X_W requires two embedded agents,
 1448 as shown in Fig. 11. The first agent (Fig. 11a) selects a recorded case based on the current
 1449 percept, and sequentially copies the actions specified by that case into X_W . The “world” of
 1450 this agent consists of X_P , M_{case} and X_W ; its “perception” kernel selects the case from M_{case}
 1451 for which the initial state is closest to the current percept x_P , its “decision” kernel selects
 1452 records from this case in sequence and its “action” kernel writes the action $g(t_F)$ specified
 1453 by the selected case into X_W . The process executed by this agent requires a time step,
 1454 i.e. one increment of t . The second agent (Fig. 11b) has a switching function analogous
 1455 to the attention-switching dyad in Fig. 10a: it compares the current percept $x_P(t)$ to the
 1456 currently-selected case record, setting $c = 1$ when the case is initially selected and setting
 1457 $c = 0$ if the distance between the states of either the object or reference components of $x_P(t)$

1458 and their states as specified by the currently-selected case record exceeds some threshold.
 1459 Setting $c = 0$ in response to such an expectation violation during case execution restores
 1460 D to its usual function. Maintaining temporal synchrony requires that the overall counter
 1461 t advances only when D executes as discussed above; this requirement can be met if D is
 1462 regarded as acting instantaneously when $c = 1$ and the action g to be selected is specified
 1463 by X_W , i.e. when action is performed “automatically.” In this case interrupting execution
 1464 of a case must be regarded as requiring one time step, after which no action is selected.



1465 *Fig. 11:* a) Selection of a case and case-record for execution based on the current
 1466 percept. This action does not enable case execution. b) Enabling or disabling
 1467 case execution by setting or resetting the control bit c based on a comparison
 1468 of current and expected percepts during case execution.

1469 The processes illustrated in Fig. 11 only execute a previous case verbatim. Interrupting
 1470 execution of a case initiates a search for a new case that is a better fit to the current per-
 1471 cept $x_P(t)$. A more intelligent case-based planner can be constructed by incorporating an
 1472 additional agent capable of modifying the currently-selected case record based on $x_P(t)$ and
 1473 information about previous component responses stored in X_C . Such modification creates

1474 a new case, which is then recorded in M_{case} . A second natural extension would incorporate
1475 a “meta” agent capable of comparing multiple cases to identify shared perception-action
1476 dependencies. A case comparator of this kind is the minimal structure needed to recognize
1477 relationships between events occurring in different orders or with different numbers of in-
1478 tervening events; hence it is the minimal structure needed to implement a “temporal map”
1479 as described by Balsam and Gallistel (2009).

1480 5 Conclusion

1481 We have shown three things in this paper. First, the CA formalism introduced by Hoffman
1482 and Prakash (2014) is both powerful and non-trivial. Even “agents” comprising only a
1483 handful of bits exhibit surprisingly complex behavior. A three-bit agent can implement a
1484 Toffoli gate, so networks of three-bit agents can compute any computable function, and
1485 can even do so reversibly. More intriguing are the hints that networks of simple agents
1486 exhibit dynamical symmetries that also characterize geometry. This result comports well
1487 with current efforts by physicists to derive the familiar geometry of spacetime from the
1488 symmetries of information exchange between simple processing units (e.g. Tegmark, 2015).
1489 We are currently working toward a full description of spacetime constructed entirely within
1490 the CA framework.

1491 We have, second, shown that concept of “fitness” as connectivity emerges naturally when
1492 networks of interacting RCAs are considered. This fitness concept accords well with estab-
1493 lished concepts of centrality developed in the theory of social and other complex networks.
1494 By expressing fitness with the CA framework, we free ITP from any need to rely on an
1495 externally-stipulated fitness function. Computational experiments to characterize the con-
1496 ditions in which preferential attachment and hence high-connectivity individuals emerge in
1497 networks of interacting RCAs are being designed.

1498 Our third result is that networks of RCAs can, at least in principle, implement sophisticated
1499 cognitive processes including attention, categorization and planning. This result fleshes out
1500 the central concepts of ITP: that experience is an *interface* onto an ontologically-ambiguous
1501 world, and that “objects” and “causal relations” are patterns of positive and negative cor-
1502 relations between experiences. It highlights the critical role played by aspects of experience
1503 that do not change, and hence serve as “context” or, more formally, reference frames rel-
1504 ative to which aspects of experience that do change can be classified and analyzed. Here
1505 again, our result comports well with recent work in physics, where with the rise of quantum
1506 information theory, the roles of reference frames in defining what can and cannot be known
1507 or communicated about a physical situation have taken on new prominence (e.g. Bartlett,
1508 Rudolph and Spekkens, 2007). A substantial program of simulation development and test-
1509 ing is clearly required to evaluate, in structured and eventually in open environments, the
1510 formal models of memory, attention, categorization and planning developed here. The level
1511 of complexity at which such models can feasibly be implemented remains unclear. We hope,
1512 however, to be able to fully characterize the reference frames required to support relatively

1513 simple behaviors in relatively simple environments, and to use this information to formulate
1514 predictions testable in more complex systems.

1515 The CA framework is, as we have emphasized, a minimal formal framework for under-
1516 standing cognition and agency. While debates about the structure and content of memory
1517 - and implicitly, experience - have dominated cognitive science for decades (e.g. Gibson,
1518 1979; Fodor and Pylyshyn, 1988; Anderson, 2003), these debates have generally been con-
1519 ducted either informally or in the context of complex, conceptually open-ended modeling
1520 paradigms. Our results, together with those of Friston and colleagues using the predictive
1521 coding and adaptive inference framework, show that cognition and agency can be addressed
1522 in conceptually very simple terms. The primary task of an organism in an environment
1523 is to regulate its interactions with the environment, by behaving appropriately, in order
1524 to maintain an environmental state conducive its own homeostasis. As Conant and Ashby
1525 (1970) showed and Friston (2010; 2013) has significantly elaborated, effective regulation
1526 of the environment requires a statistically well-founded model of the environment. Consis-
1527 tency with ITP requires that such models treat the environment as open, in which case they
1528 can be at best satisficing. The results obtained here, together with those of Friston (2013)
1529 and Friston, Levin, Sengupta and Pezzulo (2015), offer an outline of how such models may
1530 be constructed in a way that is consistent with ITP, but many details remain to be worked
1531 out. A thorough treatment of both evolutionary and developmental processes from both
1532 extrinsic and intrinsic perspectives is needed to understand the kinds of worlds W in which
1533 complex networks of interdependent RCAs can be expected to appear.

1534 We have largely deferred the question of motivation. As mentioned in §4.3 above, ratio-
1535 nal agents exhibit curiosity and hence explore their environments to discover sources of
1536 “good” experiences, which in a typical W may lie very near sources of “bad” experiences.
1537 As Gottlieb, Oudeyer, Lopes and Baranes (2013) emphasize, however, rational agents do
1538 not exhibit unlimited curiosity, as this can lead to expending all available resources at-
1539 tempting to solve unsolvable problems or learn unlearnable information. Understanding
1540 and modeling motivation requires not only a formal characterization of resources and their
1541 use, but also a formal model of reward, its representation, and its roles in both extrinsic
1542 and intrinsic motivation. The distinction between the “pragmatic” and “epistemic” values
1543 of information (Friston et al., 2015) is useful here; the current framework models the effects
1544 of this distinction in terms of attention switching, but not its origin. Both developmental
1545 robotics (e.g. Cangelosi and Schlesinger, 2015) and the neuroscience of the reward system
1546 (e.g. Berridge, and Kringelbach, 2013) provide empirical avenues to pursue in this regard.

1547 We have also, and more importantly from an architectural perspective, deferred the task
1548 of constructing a full theory of RCA networks and RCA combinations. Developing such
1549 a theory will require addressing such questions as whether RCA networks can in general
1550 be considered locally hierarchical, whether the action spaces G of complex RCAs require
1551 structures, for example to represent fully automatized action patterns, analogous to the
1552 structures in X described here, and how to explicitly define D kernels in complex RCAs. It
1553 will also require understanding how the time counters (i.e. t parameters) of complex RCAs
1554 relate to those of their component RCAs, a question that has been elided here by assuming

1555 that all processes “inside” X are synchronous. Answering such questions may well depend
1556 on resolving at least some of the issues having to do with fitness and motivation mentioned
1557 above. We expect, however, that their answers will shed light on such questions as whether
1558 complex RCAs can in some cases be regarded as unaware of the experiences - e.g. the
1559 percepts or memories - of their component RCAs and how the actions of complex RCAs
1560 depend, or not, on the actions of their component RCAs.

1561 As CAs and hence RCAs are intended, from the outset, to represent *conscious* agents, it
1562 is natural to ask what the behavior of networks of RCAs can tell us about consciousness.
1563 Here two results stand out. The first is that an agent cannot, without violating ITP,
1564 distinguish the world outside of her experience from another conscious agent. While this
1565 follows from the ontological principle of conscious realism of Hoffman and Prakash (2014),
1566 it equally follows from the impossibility, within ITP, of determining that the “world” has
1567 non-Markovian dynamics. The second is that agents can be expected to be aware of more
1568 than they can report. This seems paradoxical if awareness is equated with reportability,
1569 but makes sense when the attentional resources that would be required to enable reporting
1570 of all experiences are taken into account.

1571 While examining specific cases of successful and unsuccessful behavior in well-defined worlds
1572 requires addressing the issues of motivation and multi-agent combination highlighted above,
1573 two substantial conceptual issues stand out. The first is that the CA formalism, in contrast
1574 to either standard neural network approaches or purely-functional cognitive modelling ap-
1575 proaches, enforces by its structure a focus on what a constructed agent is being modelled
1576 as experiencing. The CA formalism itself requires that the decision kernel D acts on the
1577 space of experiences X ; hence whatever D acts on must be in X and therefore must be an
1578 experience. Constructing complex memory structures in X in order to make them available
1579 to D is, given this constraint, proposing the hypothesis that the contents of such struc-
1580 tures are experienced. Experienced by whom? Here the second issue becomes relevant.
1581 As discussed in §3.2, discussions of consciousness have often assumed, explicitly or more
1582 typically implicitly, that “low-level” experiences combine in some straightforward way into
1583 “higher-level” experiences. The phenomenal unity of ordinary, waking human experience
1584 is assumed by many to indicate that there is only one relevant “level” of experience, the
1585 level of the whole organism (or often, just its brain). With this assumption, proper com-
1586 ponents of the human neurocognitive system cannot themselves be experiencers; that this
1587 is the case is treated as axiomatic, for example, in Integrated Information Theory (Tononi
1588 and Koch, 2015; see Cerullo, 2015 for a critique of this assumption in the IIT context).
1589 If complex experiencers are networks of RCAs, however, this assumption cannot be cor-
1590 rect: all RCAs, even the simplest ones, experience *something*. If complex experiencers are
1591 networks of RCAs, there is also no reason to assume that “higher-level” experiences are in
1592 any straightforward sense combinations of “lower-level” ones. Unless RCA combinations are
1593 simple Cartesian products, high-level experiences will in general not be uniquely predictable
1594 from low-level experiences or vice-versa. If complex experiencers are only approximately
1595 hierarchical rich-club networks of RCAs, the assumption that experiences should in general
1596 be straightforwardly combinatoric is almost certainly wrong.

1597 That said, it is worth re-emphasizing that the CA framework is not, and is not intended to
1598 be, a theory of consciousness *per se*. The CA framework says nothing about the *nature* of
1599 experience. It says nothing about qualia; it simply assumes that qualia exist, that agents
1600 experience them, and that they can be tokened by elements of X . The CA framework is,
1601 instead, a formal framework for modelling conscious agents and their interactions that en-
1602 forces consistency with ITP. By itself, the CA framework is ontologically neutral, as is ITP.
1603 When equipped with the ontological assumption of conscious realism, the CA framework
1604 becomes at least *prima facie* consistent with ontological theories that take consciousness to
1605 be an irreducible primitive. The role of the CA framework in expressing the assumptions
1606 or results of such theories can be expected to depend on the details of their ontological
1607 assumptions. Whether the CA framework fully captures the ontological assumptions of
1608 existing theories that take consciousness to be fundamental, e.g. that of Faggin (2015),
1609 remains to be determined.

1610 In summary, the CA framework, and RCA networks in particular, provide both a highly-
1611 constrained formal technology for representing cognition and a way of thinking about cogni-
1612 tion that emphasizes experience and decisions based on experience. It directly implements
1613 the ontological neutrality regarding the external world that is required by ITP. As results
1614 from physics and other disciplines render naïve or even critical realism about perceived
1615 objects and causal relations increasingly hard to sustain, this ability to model experience
1616 and decision making with no supporting ontology will become increasingly critical for psy-
1617 chology and for the biosciences in general.

1618 Acknowledgements

1619 The authors thank Federico Faggin and Robert Prentner for discussions of the ideas in this
1620 paper and The Federico and Elvia Faggin Foundation for financial support. Thanks also
1621 to the reviewers for their constructive comments.

1622 References

- 1623 Adams, R. A., Friston, K. J. and Bastos, A. M. (2015). Active inference, predictive coding
1624 and cortical architecture. In M. F. Casanova and I. Opris (Eds) *Recent Advances in the*
1625 *Modular Organization of the Cortex*. Berlin: Springer (pp. 97-121).
- 1626 Adolphs, R. (2003). Cognitive neuroscience of human social behavior. *Nature Reviews*
1627 *Neuroscience* 4, 165-178.
- 1628 Adolphs, R. (2009). The social brain: Neural basis for social knowledge. *Annual Review of*
1629 *Psychology* 60, 693-716.
- 1630 Agrawal, H. (2002). Extreme self-organization in networks constructed from gene expression
1631 data. *Physical Review Letters* 89, 268702.

- 1632 Anderson, M. L. (2003). Embodied cognition: A field guide. *Artificial Intelligence* 149,
1633 91-130.
- 1634 Ashby, W. R. (1956). *Introduction to Cybernetics*. London: Chapman and Hall.
- 1635 Aspect, A., Dalibard, J. and Roger, G. (1982). Experimental test of Bell's inequalities using
1636 time-varying analyzers. *Physical Review Letters* 49, 1804-1807.
- 1637 Baars, B. J., Franklin, S. and Ramsoy, T. Z. (2013). Global workspace dynamics: Cortical
1638 "binding and propagation" enables conscious contents. *Frontiers in Psychology* 4, Article
1639 # 200.
- 1640 Baillargeon, R. (2008). Innate ideas revisited: For a principle of persistence in infants
1641 physical reasoning. *Perspectives on Psychological Science* 3, 2-13.
- 1642 Barabási, A.-L. (2009). Scale-free networks: A decade and beyond. *Science* 325, 412-413.
- 1643 Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *Science*
1644 286, 509-512.
- 1645 Barabási, A.-L. and Oltvai, Z. N. (2004). Network biology: Understanding the cell's func-
1646 tional organization. *Nature Reviews Genetics* 5, 101-114.
- 1647 Bartlett, S. D., Rudolph, T. and Spekkens, R. W. (2007). Reference frames, superselection
1648 rules, and quantum information. *Reviews of Modern Physics* 79, 555-609.
- 1649 Bassett, D. S. and Bullmore, E. (2006). Small world brain networks. *The Neuroscientist*
1650 12, 512-523.
- 1651 Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P. and Friston, K. J.
1652 (2012). Canonical microcircuits for predictive coding. *Neuron* 76, 695-711.
- 1653 Bennett, B. M., Hoffman, D. D. and Prakash, C. (1989). *Observer Mechanics: A Formal*
1654 *Theory of Perception*. Academic Press.
- 1655 Berridge, K. C. and Kringelbach, M. L. (2013). Neuroscience of affect: brain mechanisms
1656 of pleasure and displeasure. *Current Opinion in Neurobiology* 23, 294-303.
- 1657 Boly, M., Sanders, R. D., Mashour, G. A. and Laureys, S. (2013). Consciousness and
1658 responsiveness: Lessons from anaesthesia and the vegetative state. *Current Opinion in*
1659 *Anesthesiology* 26, 444-449.
- 1660 Börner, K., Sanyal, S. and Vespignani, A. (2007). Network science. *Annual Review of*
1661 *Information Science and Technology* 41, 537-607.
- 1662 Bruza, P. D., Kitto, K., Ramm, B. J. and Sitbon, L. (2015). A probabilistic framework
1663 for analysing the compositionality of conceptual combinations. *Journal of Mathematical*
1664 *Psychology* 67, 26-38.
- 1665 Burgess, P. W. and Wu, H-C. (2013). Rostral prefrontal cortex (Brodmann area 10):
1666 Metacognition in the brain. In D. T. Stuss and R. T. Knight (Eds.) *Principles of Frontal*
1667 *Lobe Function, 2nd Ed.* New York: Oxford University Press (pp. 524-534).

- 1668 Cangelosi, A. and Schlesinger, M. (2015). *Developmental Robotics: From Babies to Robots*.
1669 Cambridge, MA: MIT Press.
- 1670 Cannon, T. D. (2015). How schizophrenia develops: Cognitive and brain mechanisms
1671 underlying onset of psychosis. *Trends in Cognitive Science* 19, 744-756.
- 1672 Cerullo, M. A. (2015). The problem with Phi: A critique of Integrated Information Theory.
1673 *PLoS Computational Biology* 11, e1004286.
- 1674 Colizza, V., Flammini, A., Serrano, M. A. and Vespignani, A. (2006). Detecting rich-club
1675 ordering in complex networks. *Nature Physics* 2, 110-115.
- 1676 Conant, R. C. and Ashby, W. R. (1970). Every good regulator of a system must be a model
1677 of that system. *International Journal of Systems Science* 1, 89-97.
- 1678 Conway, J. and Kochen, S. (2006). The free will theorem. *Foundations of Physics* 36,
1679 1441-1473.
- 1680 Craig, A. D. (2010). The sentient self. *Brain Structure and Function* 214, 563-577.
- 1681 Cummins, R. (1977). Programs in the explanation of behavior. *Philosophy of Science* 44,
1682 269-287.
- 1683 Damasio, A. (1999). *The Feeling of What Happens: Body and Emotion in the Making of*
1684 *Consciousness*. Orlando, FL: Harcourt.
- 1685 Debruyne, H, Portzky, M., Van den Eynde, F. and Audenaert, K. (2009). Cotards syn-
1686 drome: A review. *Current Psychiatry Reports* 11, 197-202.
- 1687 Dehaene, S., Charles, L., King, J.-R. and Marti, S. (2014). Toward a computational theory
1688 of conscious processing. *Current Opinion in Neurobiology* 25, 76-84.
- 1689 Diestel, R. (2010). *Graph theory* (4th ed.). Berlin: Springer.
- 1690 Dunbar, R. I. M. (2003). The social brain: Mind, language and society in evolutionary
1691 perspective. *Annual Review of Anthropology* 32, 163-181.
- 1692 Dunbar, R. I. M., and Shultz, S. (2007). Evolution in the social brain. *Science* 317, 1344-
1693 1347.
- 1694 Eibenberger, S., Gerlich, S., Arndt, M., Mayor, M. and Txen, J. (2013). Matter-wave
1695 interference of particles selected from a molecular library with masses exceeding 10,000
1696 amu. *Physical Chemistry and Chemical Physics* 15, 14696-14700.
- 1697 Faggin, F. (2015). The nature of reality. *Atti e Memorie dell'Accademia Galileiana di*
1698 *Scienze, Lettere ed Arti*, Volume CXXVII (2014-2015). Padova: Accademia Galileiana di
1699 *Scienze, Lettere ed Arti*.
- 1700 Fields, C. (2002). Why do we talk to ourselves? *Journal of Experimental & Theoretical*
1701 *Artificial Intelligence* 14, 255-272.
- 1702 Fields, C. (2012). The very same thing: extending the object token concept to incorporate
1703 causal constraints on individual identity. *Advances in Cognitive Psychology* 8, 234-247.

- 1704 Fields, C. (2013a). A whole box of Pandoras: Systems, boundaries and free will in quantum
1705 theory. *Journal of Experimental & Theoretical Artificial Intelligence* 25, 291-302.
- 1706 Fields, C. (2013b). How humans solve the frame problem. *Journal of Experimental &*
1707 *Theoretical Artificial Intelligence* 25, 441-456.
- 1708 Fields, C. (2014). Motion, identity and the bias toward agency. *Frontiers in Human*
1709 *Neuroscience* 8, Article # 597.
- 1710 Fields, C. (2016). Building the observer into the system: Toward a realistic description of
1711 human interaction with the world. *Systems* 4, Article # 32.
- 1712 Fodor, J. A. and Pylyshyn, Z. W. (1988). Connectionism and cognitive architecture: A
1713 critical analysis. *Cognition* 28, 3-71.
- 1714 Ford, K. M. and Pylyshyn, Z. W. (Eds.) (1996). *The Robot's Dilemma Revisited*. Norwood,
1715 NJ: Ablex.
- 1716 Friston, K. (2010). The free-energy principle: A unified brain theory? *Nature Reviews*
1717 *Neuroscience* 11, 127-138.
- 1718 Friston, K. (2013). Life as we know it. *Journal of the Royal Society: Interface* 10, 20130475.
- 1719 Friston, K., Thornton, C. and Clark, A. (2012). Free-energy minimization and the dark-
1720 room problem. *Frontiers in Psychology* 3, article # 130.
- 1721 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P, O'Doherty, J. and Pezzulo, G.
1722 (2016). Active inference and learning. *Neuroscience and Biobehavioral Reviews* 68, 862-879.
- 1723 Friston, K., Levin, M., Sengupta, B. and Pezzulo, G. (2015). Knowing ones place: A
1724 free-energy approach to pattern regulation. *Journal of the Royal Society: Interface* 12,
1725 20141383.
- 1726 Friston, K., Rigoli, F., Ognibene, D., Mathys, C., FitzGerald, T. and Pezzulo, G. (2015).
1727 Active inference and epistemic value. *Cognitive Neuroscience* 6, 187-214.
- 1728 Fuchs, C. A. and Schack, R. (2013). Quantum-Bayesian coherence. *Reviews of Modern*
1729 *Physics* 85, 1693-1715.
- 1730 Gao, T., McCarthy, G. and Scholl, B. J. (2010). The wolfpack effect: Perception of animacy
1731 irresistibly influences interactive behavior. *Psychological Science* 21, 1845-1853.
- 1732 Gibson, J. J. (1979). *The Ecological Approach to Visual Perception*. Boston: Houghton-
1733 Mifflin.
- 1734 Geisler, W.S. and Diehl, R.L. (2003). A Bayesian approach to the evolution of perceptual
1735 and cognitive systems. *Cognitive Science* 27, 379-402.
- 1736 Giustina, M., Versteegh, M. A. M., Wengerowsky, S. et al. (2015). A significant-loophole-
1737 free test of Bells theorem with entangled photons. *Physical Review Letters* 115, 250401.
- 1738 Goldberg, R. P. (1974). A survey of virtual machine research. *IEEE Computer* 7(6), 34-45.

- 1739 Gottlieb, J., Oudeyer, P.-Y., Lopes, L. and Baranes, A. (2013). Information-seeking, curios-
1740 ity, and attention: Computational and neural mechanisms. *Trends in Cognitive Sciences*
1741 17, 585-593.
- 1742 Graziano, M. S. A. (2014). Speculations of the evolution of awareness. *Journal of Cognitive*
1743 *Neuroscience* 26, 1300-1304.
- 1744 He, X, Feldman, J. and Singh, M. (2015). Structure from motion without projective con-
1745 sistency. *Journal of Vision* 15, 725.
- 1746 Heavey, C. L. and Hurlburt, R. T. (2008). The phenomena of inner experience. *Conscious-*
1747 *ness and Cognition* 17, 798-810.
- 1748 Hensen, B., Bernien, H., Drau, A. E. et al. (2015). Loophole-free Bell inequality violation
1749 using electron spins separated by 1.3 kilometres. *Nature* 526, 682-686.
- 1750 Hoffman, D. D. (2016). The interface theory of perception. *Current Directions in Psycho-*
1751 *logical Science* 25, 157-161.
- 1752 Hoffman, D. D. and Prakash, C. (2014) Objects of consciousness. *Frontiers in Psychology*
1753 5, Article # 577.
- 1754 Hoffman, D. D. and Singh, M. (2012). Computational evolutionary perception. *Perception*
1755 41, 1073-1091.
- 1756 Hoffman, D. D., Singh, M. and Prakash, C. (2015). The interface theory of perception.
1757 *Psychonomic Bulletin & Review* 22, 1480-1506.
- 1758 Hohwy, J. (2016). The self-evidencing brain. *Noûs* 50, 259-285.
- 1759 Hommel, B. (2004). Event files: Feature binding in and across perception and action.
1760 *Trends in Cognitive Sciences* 8, 494-500.
- 1761 Jacques, V., Wu, E., Grosshans, F., Treussart, F., Grangier, P., Aspect, A. and Roch,
1762 J.-F. (2007). Experimental realization of Wheeler's delayed-choice gedanken experiment.
1763 *Science* 315, 966-968.
- 1764 Jennings, D. and Leifer, M. (2015). No return to classical reality. *Contemporary Physics*
1765 in press (arxiv:1501.03202).
- 1766 Kanai, R., komura, Y., Shipp, S. and Friston, K. (2015). Cerebral hierarchies: Predictive
1767 processing, precision and the pulvinar. *Philosophical Transactions of the Royal Society B*
1768 370, 20140169.
- 1769 Keifer, M. and Pulvermüller, F. (2012). Conceptual representations in mind and brain:
1770 Theoretical developments, current evidence and future directions. *Cortex* 7, 805-825.
- 1771 Kitto, K. (2014). A contextualised general systems theory. *Systems* 2, 541-565.
- 1772 Klein, S. B. (2014). Sameness and the self: Philosophical and psychological considerations.
1773 *Frontiers in Psychology* 5, Article # 29.
- 1774 Kolodner, J. (1993). *Case-Based Reasoning*. San Mateo, CA: Morgan Kaufmann.

- 1775 Koenderink, J. J. (2014). The all seeing eye? *Perception* 43, 1-6.
- 1776 Koenderink, J. J., van Doorn, A. J. and Todd, J. T. (2009). Wide distribution of external
1777 local sign in the normal population. *Psychological Research* 73, 14-22.
- 1778 Landauer, R. (1961). Irreversibility and heat generation in the computing process. *IBM J.*
1779 *Research Development* 5, 183-195.
- 1780 Landauer, R. (1999). Information is a physical entity. *Physica A* 263, 63-67.
- 1781 Landsman, N. P. (2007). Between classical and quantum. In: J. Butterfield and J. Earman
1782 (Eds) *Handbook of the Philosophy of Science: Philosophy of Physics*. Amsterdam: Elsevier.
1783 pp 417553.
- 1784 Lloyd, S. (2012). A Turing test for free will. *Philosophical Transactions of the Royal Society*
1785 *A* 370, 3597-3610.
- 1786 Maloney, L. T. and Zhang, H. (2010). Decision-theoretic models of visual perception and
1787 action. *Vision Research* 50, 2362-2374.
- 1788 Manning, A. G., Khakimov, R. I., Dall, R. G. and Truscott, A. G. (2015). Wheelers
1789 delayed-choice gedanken experiment with a single atom. *Nature Physics* 11, 539-542.
- 1790 Mark, J. T., Marion, B. B., and Hoffman, D. D. (2010). Natural selection and veridical
1791 perceptions. *Journal of Theoretical Biology* 266, 504-515.
- 1792 Marr, D. (1982). *Vision*. San Francisco: Freeman.
- 1793 Martin, A. (2007). The representation of object concepts in the brain. *Annual Review of*
1794 *Psychology* 58, 25-45.
- 1795 McCarthy, J. and Hayes, P. J. (1969). Some philosophical problems from the standpoint
1796 of artificial intelligence. In D. Michie and B. Meltzer (Eds.) *Machine intelligence, Vol. 4*.
1797 Edinburgh: Edinburgh University Press (pp. 463-502).
- 1798 Mermin, N. D. (1985). Is the moon there when nobody looks? Reality and the quantum
1799 theory. *Physics Today* 38(4), 38-47.
- 1800 Merton, R. K. (1968). The Matthew effect in science. *Science* 159, 56-63.
- 1801 Metzinger, T. (2011). The no-self alternative. In Gallagher, S (Ed.) *The Oxford Handbook*
1802 *of the Self*. Oxford: Oxford University Press (pp 287-305).
- 1803 Newman, M. E. J. (2001). The structure of scientific collaboration networks. *Proceedings*
1804 *of the National Academy of Sciences USA* 98, 404-409.
- 1805 Overgaard, M. (2011). Visual experience and blindsight: A methodological review. *Exper-*
1806 *imental Brain Research* 209, 473-479.
- 1807 Palmer, S. E. (1999). *Vision Science: Photons to Phenomenology*. Cambridge, MA: MIT
1808 Press.
- 1809 Parthasarathy, K. R. (2005). *Introduction to Probability and Measure*. Gurgaon, India:
1810 Hindustan Book Agency.

- 1811 Pattee, H. H. (2001). The physics of symbols: Bridging the epistemic cut. *Biosystems* 60,
1812 5-21.
- 1813 Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible*
1814 *Inference*. San Mateo CA: Morgan Kaufmann.
- 1815 Peil, K. (2015). Emotional sentience and the nature of phenomenal experience. *Progress*
1816 *in Biophysics and Molecular Biology* 119, 545-562.
- 1817 Pizlo, Z., Li, Y., Sawada, T. and Steinman, R.M. (2014). *Making a Machine that Sees Like*
1818 *Us*. New York: Oxford University Press.
- 1819 Polanyi, M. (1968). Lifes irreducible structure. *Science* 160, 1308-1312.
- 1820 Pont, S. C., Nefs, H. T., van doorn, A. J., Wijntjes, M. W. A., te Pas, S. F., de Ridder, H.
1821 and Koenderink, J. J. (2012). *Seeing and Perceiving* 25, 339-349.
- 1822 Popper, K. (1963). *Conjectures and Refutations: The Growth of Scientific Knowledge*.
1823 London: Routledge & Kegan Paul.
- 1824 Pothos, E. M. and Busemeyer, J. M. (2013). Can quantum probability provide a new
1825 direction for cognitive modeling? *Behavioral and Brain Sciences* 36, 255-327.
- 1826 Prakash, C. and Hoffman, D. D. (2016). Structure invention by conscious agents. In review.
- 1827 Prakash, C., Hoffman, D. D., Stephens, K. D., Singh, M. and Fields, C. (2016). Fitness
1828 beats truth in the evolution of perception. In review.
- 1829 Pylyshyn, Z. W. (Ed.) (1987). *The Robot's Dilemma*. Norwood, NJ: Ablex.
- 1830 Riesbeck, C. K. and Schank, R. C. (1989). *Inside Case-Based Reasoning*. Hillsdale, NJ:
1831 Erlbaum.
- 1832 Rochat, P. (2012). Primordial sense of embodied self-unity. In: V. Slaughter and C. A.
1833 Brownell (Eds), *Early Development of Body Representations*. Cambridge, UK: Cambridge
1834 University Press (pp. 3-18).
- 1835 Rosen, R. (1986). On information and complexity. In J. L. Casti and A. Karlqvist (Eds.),
1836 *Complexity, Language, and Life: Mathematical Approaches*. Berlin: Springer (pp. 174-
1837 196).
- 1838 Rubino, G., Rozema, L. A., Feix, A., Araújo, M., Zeuner, J. M., Procopio, L. M., Brukner,
1839 Č. and Walter, P. (2016). Experimental verification of an indefinite causal order. Preprint
1840 arxiv:1608.01683v2 [quant-ph].
- 1841 Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*
1842 61, 85-117.
- 1843 Scholl, B. J. (2007). Object persistence in philosophy and psychology. *Mind and Language*
1844 22, 563-591.
- 1845 Scholl, B. J. and Gao, T. (2013). Perceiving animacy and intentionality: Visual processing
1846 or higher-level judgment?. In M. D. Rutherford and V. A. Kuhlmeier (Eds.) *Social Per-*
1847 *ception: Detection and Interpretation of Animacy, Agency and Intention*. Cambridge, MA:
1848 MIT Press (pp. 197-230).

- 1849 Scholl, B. J., and Tremoulet, P. (2000). Perceptual causality and animacy. *Trends in*
1850 *Cognitive Science* 4, 299309.
- 1851 Seth, A. K. (2013). Interoceptive inference, emotion, and the embodied self. *Trends in*
1852 *Cognitive Sciences* 17, 565-573.
- 1853 Shalm, L. K., Meyer-Scott, E., Christensen, B. G. et al. (2015). A strong loophole-free test
1854 of local realism. *Physical Review Letters* 115, 250402.
- 1855 Shipp, S., Adams, R. A. and Friston, K. J. (2013). Reflections on agranular architecture:
1856 Predictive coding in the motor cortex. *Trends in Neuroscience* 36, 706-716.
- 1857 Simons, D. J. and Ambinder, M. S. (2005). Change blindness: Theory and consequences.
1858 *Current Directions in Psychological Science* 14(1), 44-48.
- 1859 Simons, J. S., Henson, R. N. A., Gilbert, S. J. and Fletcher, P. C. (2008). Separable
1860 forms of reality monitoring supported by anterior prefrontal cortex. *Journal of Cognitive*
1861 *Neuroscience* 20, 447-457.
- 1862 Smith, J. E. and Nair, R. (2005). The architecture of virtual machines. *IEEE Computer*
1863 38(5), 32-38.
- 1864 Steptoe, A., Shankar, A., Demakakos, P. and Wardle, J. (2013). Social isolation, loneliness,
1865 and all-cause mortality in older men and women. *Proceedings of the National Academy of*
1866 *Sciences USA* 110, 5797-5801.
- 1867 Tanenbaum, A. S. (1976). *Structured Computer Organization*. Upper Saddle River, NJ:
1868 Prentice Hall.
- 1869 Tegmark, M. (2015). Consciousness as a state of matter. *Chaos, Solitons & Fractals* 76,
1870 238-270.
- 1871 Toffoli, T. (1980). Reversible computing. In: J. W. de Bakker and J. van Leeuwen (Eds)
1872 *Automata, Languages and Programming: Lecture Notes in Computer Science, Vol. 85*.
1873 Berlin: Springer. pp. 632644.
- 1874 Tononi, G. and Koch, C. (2015). Consciousness: Here, there and everywhere? *Philosophical*
1875 *Transactions of the Royal Society B* 370, 20140167.
- 1876 Trivers, R. L. (2011). *The Folly of Fools*. New York: Basic Books.
- 1877 Turing, A. R. (1936). On computable numbers, with an application to the Entschei-
1878 dungsproblem. *Proceedings of the London Mathematical Society, Series 2* 442, 230-265.
- 1879 van den Heuvel, M. P. and Sporns, O. (2011). Rich-club organization of the human con-
1880 nectome. *Journal of Neuroscience* 31, 15775-15786.
- 1881 Vogel, E. K., Woodman, G. F. and Luck, S. J. (2006). The time course of consolidation
1882 in visual working memory. *Journal of Experimental Psychology: Human Perception and*
1883 *Performance* 32, 1436-1451.
- 1884 von Hofsten, C. (2007). Action in development. *Developmental Science* 10, 54-60.

- 1885 von Uexküll, J. (1957). A stroll through the worlds of animals and men. In: C. Schiller
1886 (Ed.) *Instinctive Behavior*. New York: van Nostrand Reinhold (pp. 5-80). Also published
1887 in *Semiotica* 89 (1992) 319-391.
- 1888 Vossel, S., Geng, J. J. and Fink, G. R. (2014). Dorsal and ventral attention systems :
1889 Distinct Neural Circuits but collaborative roles. *The Neuroscientist* 20, 150-159.
- 1890 Wang, Q., Schoenlein, R. W., Peteanu, L. A., Mathies, R. A. and Shank, C. V. (1994).
1891 Vibrationally coherent photochemistry in the femtosecond primary event of vision. *Science*
1892 266, 422-424.
- 1893 Watson, T. L., Robbins, R. A. and Best, C. T. (2014). Infant perceptual development
1894 for faces and spoken words: An integrated approach. *Developmental Psychobiology* 56,
1895 1454-1481.
- 1896 Watts, D. J and Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks.
1897 *Nature* 393, 440-442.
- 1898 Wiseman, H. (2015). Quantum physics: Death by experiment for local realism. *Nature*
1899 526, 649-650.