

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 .*

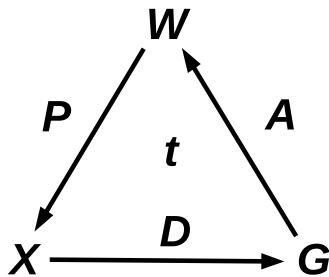
170 Less formally, if K is a Markovian kernel from B to C , then for any measurable $D \subset C$, the
171 function enacted by $x \mapsto K(x, D) \in [0, 1]$ assigns to each x in B a probability distribution

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 Let $\langle W, \mathcal{W} \rangle$, $\langle X, \mathcal{X} \rangle$ and $\langle G, \mathcal{G} \rangle$ be measurable spaces. Hoffman and Prakash (2014)
 175 defined a CA, given the measurable space $\langle W, \mathcal{W} \rangle$, as a 6-tuple $[(X, \mathcal{X}), (G, \mathcal{G}), P, D, A, t]$
 176 where $P : W \times \mathcal{X} \rightarrow [0, 1]$, $D : X \times \mathcal{G} \rightarrow [0, 1]$ and $A : G \times \mathcal{W} \rightarrow [0, 1]$ are Markovian
 177 kernels and t is a positive integer parameter. Here we explicitly include $\langle W, \mathcal{W} \rangle$ in the
 178 definition of a CA. Following Hoffman, Singh and Prakash (2015) and Prakash and Hoffman
 179 (in review), we also explicitly allow the P , D , and A kernels to depend on the elements
 180 of their respective target sets. Informally, for $x \in X$ and $g \in G$, for example, and any
 181 measurable $H \subset G$, the function enacted by $(x, g) \mapsto K(x, g, H)$ is real-valued and can
 182 be considered to be the regular conditional probability distribution $\text{Prob}(H|x, g)$ under
 183 appropriate conditions on the spaces involved (Parthasarathy, 2005). We have:

184 **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
 185 Markovian kernel $P : W \times X \rightarrow X$, D be a Markovian kernel $D : X \times G \rightarrow G$, and
 186 A be a Markovian kernel $A : G \times W \rightarrow W$. A **conscious agent (CA)** is a 7-tuple
 187 $[(X, \mathcal{X}), (G, \mathcal{G}), (W, \mathcal{W}), P, D, A, t]$, where t is a positive integer parameter.

188 The difference in representational power between the more general, target-set dependent
 189 kernels specified here and the original, here termed “forgetful,” kernels of Hoffman and
 190 Prakash (2014) is discussed below. We represent a CA as a labelled directed graph as
 191 shown in Fig. 1. This graph implies the development of a cyclic process, in which we can
 192 think of, e.g. the kernel $D : X \times G \rightarrow G$ as follows: for each instantiation g_0 of G in the
 193 immediately previous cycle, and the current instantiation of $x \in X$, $D(x, g_0; \cdot)$ gives the
 194 probability distribution of the $g \in G$ instantiated at the next step. The other kernels A
 195 and P are interpreted similarly.



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

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

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

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

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

259 2.2 Extrinsic and intrinsic perspectives

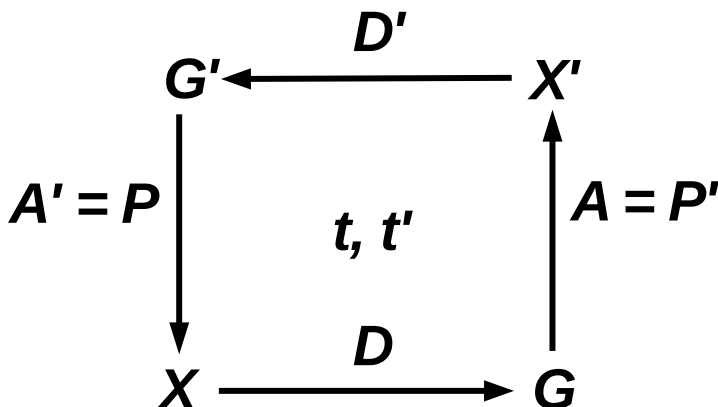
260 A central claim of ITP is that perceptual systems do not, in general, provide a veridical
 261 representation of the structure of the world; in particular, “objects” and “causal relations”
 262 appearing as experiences in X are in general not in any sense homomorphic to elements or
 263 relationships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

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

283 any experience $x \in X$; this choice is represented by the kernel D . The action A on W
 284 that the RCA is *capable* of taking is determined, in part, by the structure of W . Similarly,
 285 the action P with which W can affect the RCA is determined, in part, by the structure
 286 of the RCA. With this terminology, the central claim of ITP is that an RCA’s possible
 287 knowledge of W is completely specified by X ; the element(s) of X that are selected by P
 288 at any given t constitute the RCA’s entire experience of W at t . The structure and content
 289 of X completely specify, therefore, the intrinsic perspective of the RCA. In particular, ITP
 290 allows the RCA no independent access to the ontology of W ; consistency between intrinsic
 291 and extrinsic perspectives requires that no such access is attributed to any RCA from the
 292 latter perspective. An RCA does not, in particular, have access to the definitions of its
 293 own P , D or A kernels; hence an RCA has no way to determine whether any of them are
 294 homomorphisms. Similarly, an RCA has no access to the definitions of any other RCA’s P ,
 295 D or A kernels, or to any other RCA’s X or G . An RCA “knows” what currently appears
 296 as an experience in its own X but nothing else; as discussed in §4.1 below, for an RCA
 297 even to know what actions it has available or what actions it has taken in the past, these
 298 must be represented explicitly in X . Any structure attributed to W from the intrinsic
 299 perspective of an RCA is hypothetical in principle; such attributions of structure to W can
 300 be disconfirmed by continued observation, i.e. additional input to X , but can never be
 301 confirmed. In this sense, any RCA is in the epistemic position regarding W that Popper
 302 (1963) claims characterizes all of science.

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

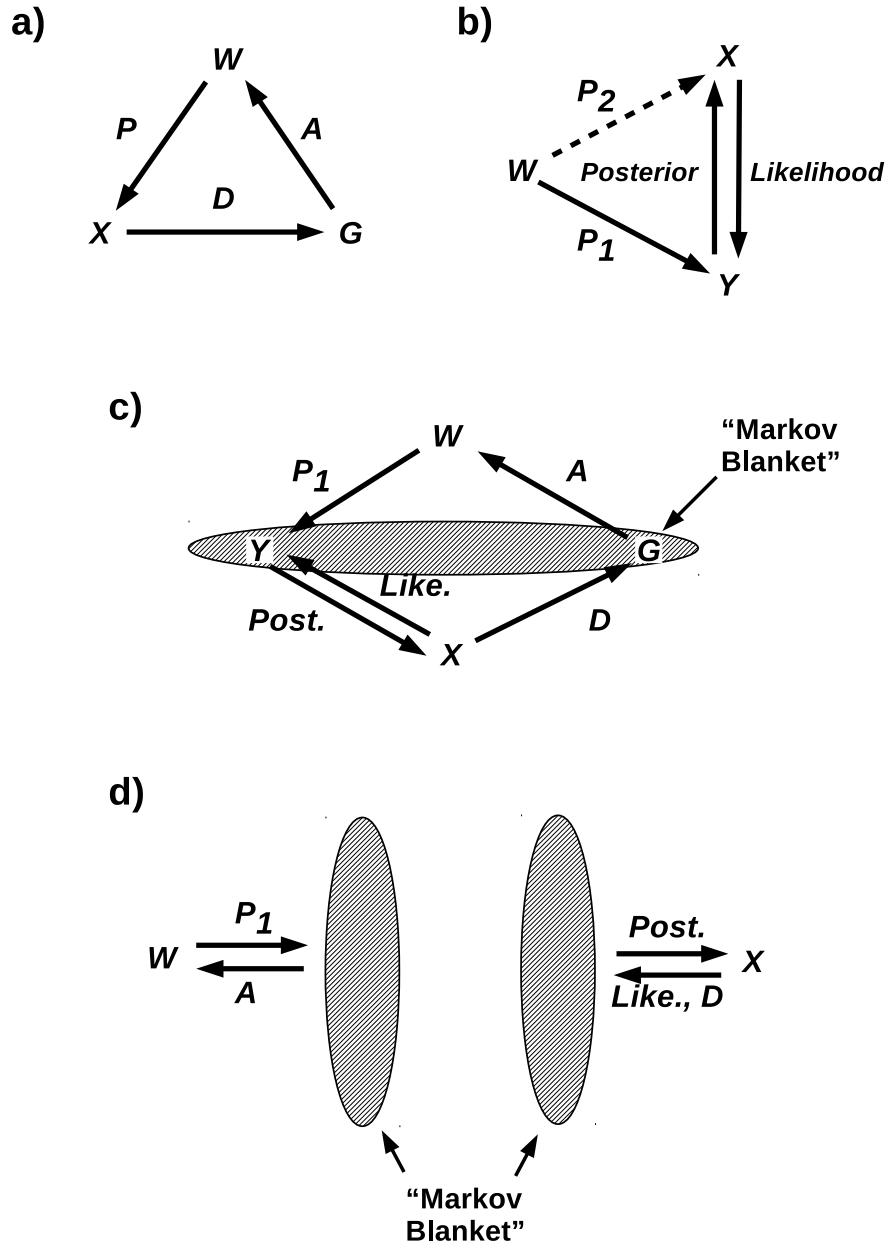
326 allowable. Every occurrence of the symbol W can, therefore, be replaced, as in Fig. 2,
 327 by an RCA. When this is done, all actions - all kernels A - act directly on the experience
 328 spaces X of other RCAs as shown in Fig. 2. If it is possible to consider any arbitrary
 329 system - any directed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



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

337 2.3 Bayesian inference and the Markov blanket

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



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

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

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

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

389 Friston (2013) argues that any random ergodic system comprising two subsystems separated
390 by a Markov blanket can be interpreted as minimizing a variational free energy that can, in
391 turn, be interpreted in Bayesian terms as a measure of expectation violation or “surprise.”
392 This Bayesian interpretation of “inference” through a Markov blanket is fully consistent
393 with the model of perceptual inference provided by the CEP framework. Conscious agents as
394 described here can, therefore, be regarded as free-energy minimizers as described by Friston
395 (2010). This formal as well as interpretational congruence between the CA framework and
396 the free-energy principle (FEP) framework of Friston (2010) is explored further below,
397 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;

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

442 **3 W from the extrinsic perspective: RCA networks** 443 **and dynamic symmetries**

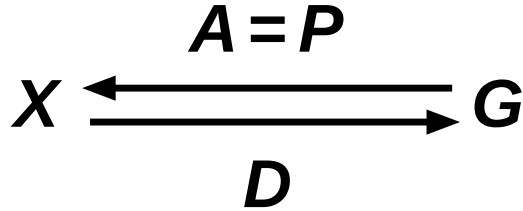
444 **3.1 Symmetric interactions**

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

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

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

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



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

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

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

488 the NOT operator,

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

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

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

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

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

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

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

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

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

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

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

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

510 It is important to emphasize that there is no sense in which the 1-bit dyad *experiences* the
 511 potential complexity of its dynamics, or in which the experience of a 1-bit dyad with one
 512 choice of D and A operators is any different from the experience of a 1-bit dyad with another
 513 choice of operators. Any 1-bit dyad has only two possible experiences, those tokened by $|0\rangle$
 514 and $|1\rangle$. The addition of memory to a CA in order to enable it to experience a *history* of
 515 states and hence relations between states from its own intrinsic perspective is discussed in
 516 §4 below.

517 The Identity and NOT operators can be expressed as “forgetful” kernels, i.e. kernels that
 518 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
 519 by Hoffman and Prakash (2014) have less representational power than the state-dependent
 520 kernels employed in the current definition of a CA. It is also worth noting that the standard
 521 AND operator taking $x(t)$ and $g(t)$ to $x(t + 1) = x(t)$ and $g(t + 1) = x(t)$ AND $g(t)$ may
 522 be represented as:
 523

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

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

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

526 The value of $G(t)$ cannot be recovered following the action of either of these operators; they
 527 are therefore logically irreversible. As each of the matrix representations of these operators
 528 has a row of all zeros, they are not Markovian. The logically irreversible, non-Markovian
 529 nature of these operators has, indeed, been a primary basis of criticisms of artificial neural
 530 network and dynamical-system models of cognition; Fodor and Pylyshyn (1988), for
 531 example, criticize such models as unable, in principle, to replicate the compositionality of
 532 Boolean operations in domains such as natural language. The standard AND operator
 533 can, however, be implemented reversibly by adding a single ancillary z bit to X , fixing its
 534 value at 0, and employing the Toffoli gate that maps $[x, y, z]$ to $[x, y, (x \text{ AND } y) \text{ XOR } z]$,
 535 where XOR is the standard exclusive OR (Toffoli, 1980). The Toffoli gate preserves the
 536 values of x and y and allows the value of z to be computed from the values of x and y ;
 537 hence it is reversible and can, therefore, be represented as a punctual Markovian kernel.
 538 The standard XOR operator employed in the Toffoli gate is equivalent to a cNOT. As any
 539 universal computing formalism must be able to compute AND, the 1-bit dynamics of Fig.
 540 4 is not computationally universal. The Toffoli gate is, however, computationally universal,

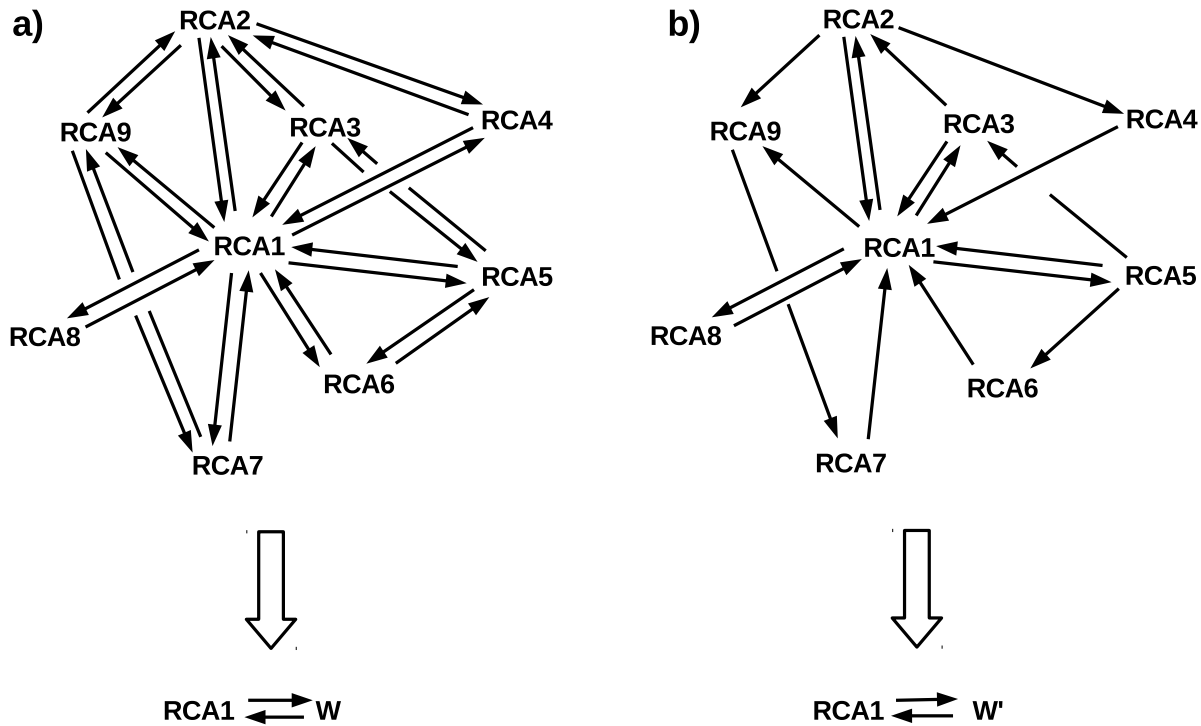
541 so adding a single ancillary bit set to 0 to each space in Fig. 4 is sufficient to achieve
542 universality.

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

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

565 3.2 Asymmetric interactions and RCA combinations

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

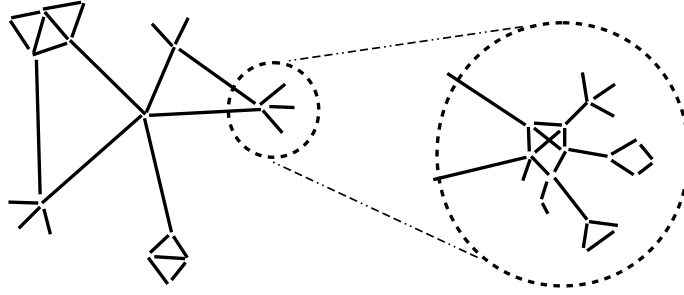


576 *Fig. 5:* a) Nine bidirectionally interacting RCAs, equivalent to a single RCA
 577 interacting with its “world” W and hence to a single CA. b) A network similar
 578 to that in a), except that some interactions are not bidirectional. Here again,
 579 the RCA network is equivalent to a single RCA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

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

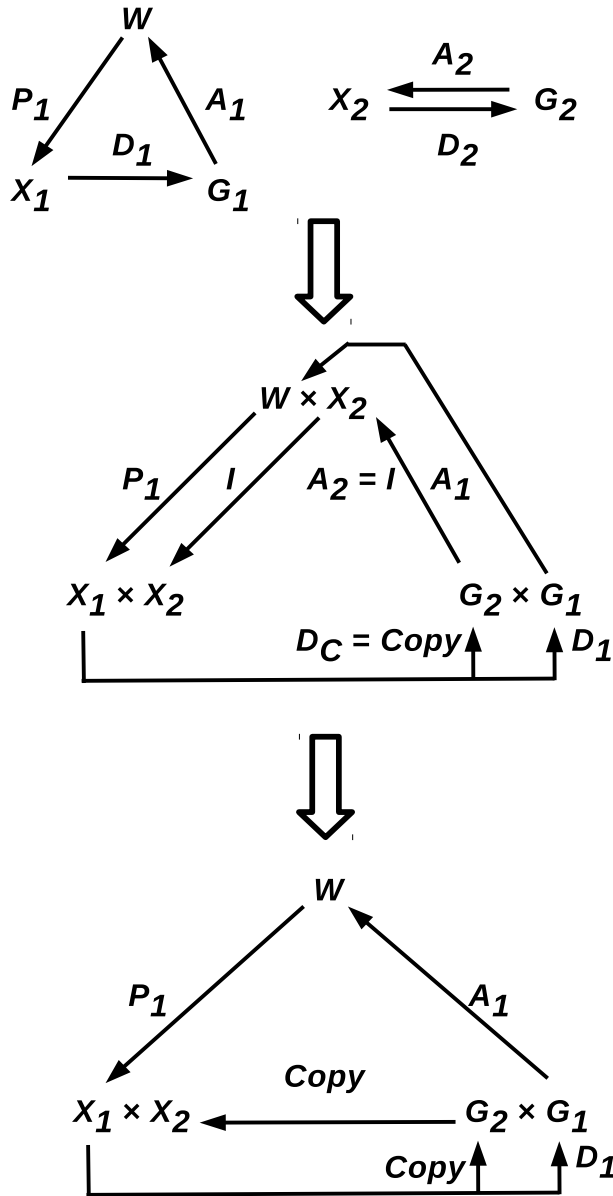
593 the “world”; indeed an arbitrarily large number of alternative structures could produce
594 the same inputs to RCA1 and hence the same sequence of experiences for RCA1. RCA1
595 cannot, in particular, determine what other RCA(s) it is interacting with at any particular
596 “time” t as measured by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

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



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

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



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

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

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

698 3.3 Connectivity and fitness

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

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

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

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

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

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

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

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

779 The identification of fitness with connectivity provides a straightforward solution to the
780 “dark room” problem faced by uncertainty-minimization systems (e.g. Friston, Thornton
781 and Clark, 2012). Dark rooms do not contain opportunities to create or maintain connec-
782 tions; therefore fitness-optimizing systems can be expected to avoid them. This solution
783 complements that of Friston, Thornton and Clark (2012), who emphasize the costs to
784 homeostasis of remaining in a dark room. Here again, interactivity and maintenance of
785 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

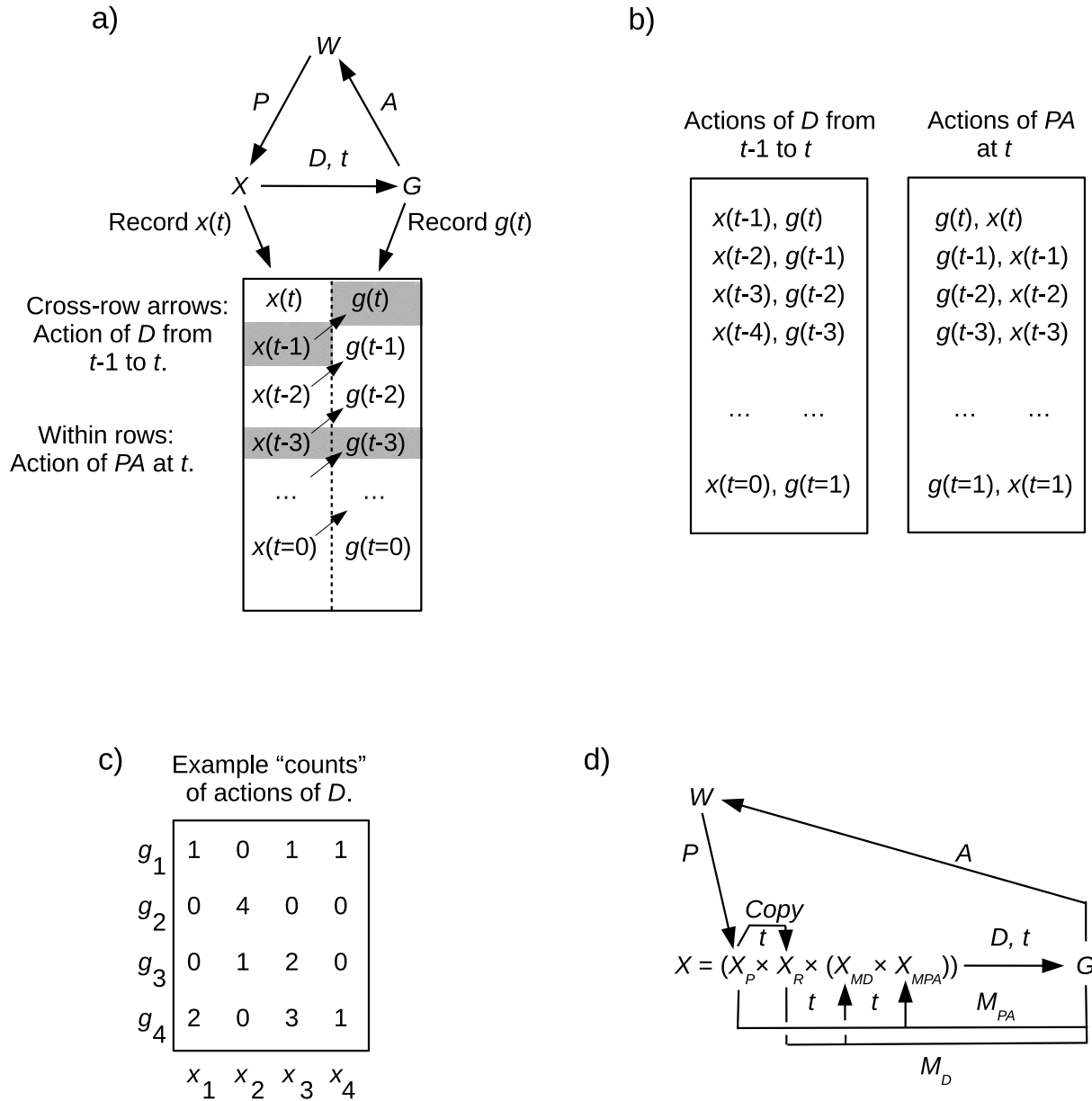
826 “collapse of the wavefunction” in quantum theory (for an accessible and thorough review,
827 see Landsman, 2007) and is often associated with the operation of free will (reviewed by
828 Fields, 2013a). We adopt this association of “collapse” with free will here: the RCA renders
829 P punctual by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

841 4.2 Memory

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

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

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



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

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

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

915 as in the middle part of Fig. 7.

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

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

942 4.3 Predictive coding, goals and active inference

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

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

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

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

996 function), that this “reality monitoring” capability appears to be highly but not exclusively
997 localized to rostral prefrontal cortex, and that disruption of this capability correlates with
998 psychosis (Simons, Gilbert, Henson and Fletcher, 2008; Burgess and Wu, 2013; Cannon,
999 2015). Humans also experientially distinguish short-term “working” memories from long-
1000 term memories. We predict that specific monitoring capabilities provide the experiential
1001 distinctions between short- (e.g. X_R) and long-term (e.g. X_{MD} and X_{MPA}) memories and
1002 distinguish functionally-distinct long-term memory components from each other. From a
1003 formal standpoint, such distinguishing tags can be considered to be additional elements in
1004 each vector in each of the derived vector spaces; while such tags play no explicit role in the
1005 processing described below, their existence will be assumed.

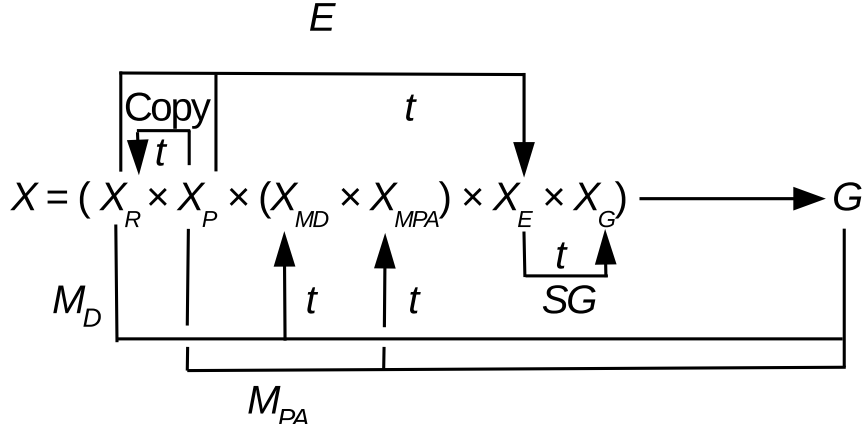
1006 As the memories X_{MD} and X_{MPA} and hence the conditional probability distributions
1007 $\text{Prob}_D(x(t), g(t)|x(t-1), g(t-1))$ and $\text{Prob}_{PA}(x(t), g(t)|x(t-1), g(t-1))$ contain informa-
1008 tion about the observer’s entire experience of the world, they enable differential responses
1009 to $x_R - g$ or $g - x_P$ pairings that evoke different degrees of “surprise” by either confirming
1010 or disconfirming previous associations to different extents. We note that the term ‘surprise’
1011 is being used here in its informal sense of an *experienced* departure from expectations, not
1012 in the technical sense employed by Friston (2010; 2013; see also Friston et al., 2015; Fris-
1013 ton et al., 2016) to refer to an event that causes or threatens to cause a departure from
1014 homeostasis and hence has negative consequences for fitness. To implement such differen-
1015 tial responses to surprise, it is natural to choose functions for updating these conditional
1016 probability distributions that depend on the vector distance(s) between the percept x_R (for
1017 $\text{Prob}_D(x(t), g(t)|x(t-1), g(t-1))$) or x_P (for $\text{Prob}_{PA}(x(t), g(t)|x(t-1), g(t-1))$) and the
1018 percept(s) previously associated, within X_{MD} and X_{MPA} respectively, with g . Functions
1019 can clearly chosen that either enhance or suppress memories of surprising events. This
1020 generalization requires no additional components or elements within X ; hence it enhances
1021 function without altering the architecture.

1022 The simplest possible action is no action: the agent merely observes the world. The extremal
1023 outcomes of such observation are on the one hand James’ “blooming, buzzing confusion,” i.e.
1024 a completely random $x_P(t)$, and on the other stasis, a fixed and invariant $x_P(t)$. Memory is
1025 obviously useless in either case; indeed, the latter corresponds to the “dark room” situation
1026 discussed above. Memory becomes useful if a world on which no action is taken generates
1027 some number of the possible percepts significantly more often than the others. The same
1028 is true in the case of any other constantly-repeated action. It is equivalent to say: any
1029 action which, when repeated indefinitely, is followed by either random or static percepts
1030 is a useless action to take. Such an action has no “epistemic value” in the sense used by
1031 Friston et al. (2015). Randomness and stasis may be useful as *components* of experience -
1032 indeed as discussed below, stasis is a *necessary* component of useful experience - but only
1033 when embedded in non-random, non-static contexts. Let us assume, therefore, that RCAs
1034 of interest are embedded in W s that generate non-random, non-static percepts in response
1035 to all actions. Note that this assumption is consistent with ITP: it does not require either
1036 P or A to respect the causal structure of W .

1037 In a non-random, non-static world, the memories X_{MD} and X_{MPA} provide a basis for

1038 predictive coding: the probability assigned to an action g at $t + 1$ can depend on the vector
 1039 difference between the current percept $x_P(t)$ and previous percepts either immediately-
 1040 antecedent or immediately-consequence to actions like g . A percept $x_P(t)$ can, in this case,
 1041 “predict” an action $g(t + 1)$ that is “expected,” on the basis of the probabilities stored
 1042 in X_{MPA} , to result in a subsequent percept $x_P(t + 1)$ that is either similar or dissimilar
 1043 to $x_P(t)$. Assigning high probabilities to actions at $t + 1$ expected to result in percepts
 1044 similar to $x_P(t)$ is implicitly “evaluating” $x_P(t)$ as in some sense “good” or “desirable,”
 1045 while assigning low probabilities to actions at $t + 1$ expected to result in percepts similar to
 1046 $x_P(t)$ is implicitly evaluating $x_P(t)$ as in some sense bad or undesirable. These operational
 1047 senses of “good” and “bad” percepts are consistent with the senses of “good” and “bad”
 1048 percepts as enhancing or threatening the maintenance of homeostasis employed by Friston
 1049 (2010; 2013). A “bad” experience in this operational sense is a outcome that an agent
 1050 did not expect to experience, i.e. a stressor such as being hungry or poor, on the basis
 1051 of the implicit “model” of W encoded by the probability distributions contained in the
 1052 memories X_{MD} and X_{MPA} . In the limit, a maximally “bad” experience is one that violates
 1053 the fundamental expectation that experiences will continue that is encoded by all non-
 1054 zero values of the subjective probabilities $\text{Prob}_D(x, g)|_{t=T}$ and $\text{Prob}_{PA}(x, g)|_{t=T}$; such an
 1055 experience destroys connectivity between the agent in question and the surrounding RCA
 1056 network (i.e. the agent’s W), setting the agent’s fitness to zero and corresponding to the
 1057 “death” of the agent as discussed in §3.3 above.

1058 This evaluative function can be made explicit by representing it as a distinct operation. To
 1059 do this, we add a further memory component X_E to X . To allow for the possibility that
 1060 an observer has “innate” biases toward or against particular percepts, we consider X_E to
 1061 comprise two probability distributions, $\text{Prob}_{good}(x_P)$ and $\text{Prob}_{bad}(x_P)$, with *a priori* values
 1062 fixed at $t = 0$. Such innate evaluation biases can be considered to be innate “preferences”
 1063 or “beliefs” as they often are in the infant-cognition literature (e.g. Baillargeon, 2008;
 1064 Watson, Robbins and Best, 2014). We represent the evaluation operation E as having two
 1065 components $E = (E_{good}, E_{bad})$, where E_{good} is a punctual kernel $X_P \times X_R \times E \rightarrow E$ that
 1066 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
 1067 updates $\text{Prob}_{bad}(x_P)$ at each t . For simplicity, we assume that E_{good} increases $\text{Prob}_{good}(x_P)$
 1068 by a factor ≥ 1 that approaches unity as $\text{Prob}_{good}(x_P) \rightarrow 1$ whenever both $\text{Prob}_{good}(x_P(t)) >$
 1069 0 and $\text{Prob}_{good}(x_R(t)) > 0$ and that E_{bad} increases $\text{Prob}_{bad}(x_P)$ by a factor with similar
 1070 behavior whenever both $\text{Prob}_{bad}(x_P(t)) > 0$ and $\text{Prob}_{bad}(x_R(t)) > 0$. This E effectively
 1071 implements the heuristic: an experience is remembered as better if it is followed by a good
 1072 experience, and remembered as worse if it is followed by a bad experience. Note that while
 1073 this heuristic is consistent with the association of “good” and “bad” with maintaining or
 1074 not maintaining either homeostasis or connectivity as discussed above, it also allows a
 1075 given x_P to be both probably good and probably bad, a not-unrealistic situation. This
 1076 additional structure on X is summarized in Fig. 9. Extending the evaluative process from
 1077 the scalar representation provided by these probabilities to a multidimensional, i.e. vector,
 1078 representation costs memory and kernel complexity but does not change the architecture.



1079 *Fig. 9:* Adding memories for evaluations of percepts (X_E) and for a current
 1080 goal (X_G) to Fig. 7d. Connections to W have been elided for clarity.

1081 Evaluating percepts implicitly evaluates the actions that are followed by those percepts;
 1082 this implicit transfer of estimated “good” or “bad” value from percepts to actions is now
 1083 implemented by D . A “rational” D , for example, would assign high probabilities to actions
 1084 g that are associated in X_{MPA} with subsequent percepts that have high valuations in X_E .
 1085 If W is such that the relative ranking of percepts by value changes only slowly with t ,
 1086 relatively highly- and lowly-ranked percepts can be considered to be positive and negative
 1087 “goals” respectively. As Friston (2010, 2013) has emphasized, goals are effectively long-term
 1088 expectations to which an uncertainty-minimizing agent attempts to match perceptions;
 1089 Friston and colleagues call acting so as to match perceptions to goals “active inference.”
 1090 Within the CA framework, the minimal functional architecture required for active inference
 1091 is that shown in Fig. 9. Here a memory component X_G holds the current goal; it is
 1092 populated by a punctual, forgetful kernel SG acting on X_E . While SG can be taken to
 1093 choose percepts of high value as goals, its specific action can be left open. Note than in this
 1094 architecture, incremental adjustments of the “world model” X_{MPA} and “self model” X_D
 1095 are made in parallel with active inference: expectations are modified to fit perceptions even
 1096 when actions are taken to modify perceptions to fit expectations. Note also that placing
 1097 the evaluation and goal memories X_E and X_G within the experience space X is predicting
 1098 that the contents of these memories are both experienced and experienced as distinct, as
 1099 they indeed are in neurotypical humans. While the specific mechanisms implementing the
 1100 experiential distinction between these memory components remains uncharacterized, the
 1101 present framework predicts that such mechanisms exist.

1102 By iteratively constructing representations of the antecedents and consequences of actions,
 1103 the kernels M_D and M_{PA} implement a simple kind of learning. The operator E similarly

1104 implements a simple form of evaluative feedback. The action choices made by D can,
1105 therefore, progressively improve with experience. It is important to emphasize that M_D ,
1106 M_{PA} , E , SG and D are all by assumption homogeneous kernels. What changes as the
1107 system learns is not the choice function D , but the contents of the data structures – the
1108 memories X_{MD} , X_{MPA} , X_E and X_G – that serve as ancillary inputs to D . The “knowledge”
1109 of an RCA with this architecture is, therefore, entirely explicit. This is marked contrast
1110 to typical neural-network models, including recent “deep learning” models (for a recent
1111 review, see Schmidhuber, 2015), in which learning is entirely implicit and the decision rules
1112 learned are notoriously hard to reverse engineer. It is worth noting that standard neural-
1113 network models have no intrinsic perspective; as emphasized earlier, it is the requirement
1114 that an RCA learns about W from its own intrinsic perspective that forces what is learned
1115 to be made explicit in a memory located in X , i.e. in a memory encoding contents that
1116 are experienced - but are not necessarily reportable - by the RCA. While the kernels M_D ,
1117 M_{PA} , E , SG , as well as others to be introduced below, that populate explicit memories
1118 can, together with the decision kernel D be considered to encode implicit memories in the
1119 current model, the assumption that all such kernels are homogeneous implies that these
1120 implicit memories are not loci of learning. The kinds of “practised skill” memories that
1121 are canonically regarded as implicit are most naturally modelled as structures, e.g. fixed
1122 or fully-automatized learned action patterns, within the action space G in the current
1123 framework; an exploration of such structures are developed within G is beyond the present
1124 scope.

1125 It is important to note that whether D is “rational” in the sense of favoring actions that re-
1126 sult in “good” outcomes, and hence the extent to which the choices favored by D “improve”
1127 with experience, is left open within the architecture. If W is such that “good” choices cor-
1128 relate with the acquisition of resources required for survival, a basic orientation or “drive”
1129 toward increasing the average subjective valuation of “good” percepts can be expected to
1130 emerge in a population of agents whenever the required resources are scarce. Friston has
1131 argued that predictability of experience is itself the primary resource that organisms seek
1132 to maximize, and that the drive to pursue and acquire external resources can be under-
1133 stood in terms of maintaining the predictability of experiences that facilitate or enhance
1134 the maintenance of physiological homeostasis (Friston, 2010; 2013; Friston, Thornton and
1135 Clark, 2012). Reducing the uncertainty of experiences from a large environment requires
1136 extensive sampling of the environment’s behaviors and hence active exploration; effective
1137 agents in a large W can, therefore, be expected to display a “curious rationality” that
1138 maintains homeostasis while devoting significant energy to active exploration and learning
1139 (reviewed by Gottlieb, Oudeyer, Lopes and Baranes, 2013). Friston et al. (2015; 2016)
1140 make a similar point: the minimization of expected surprise in the strict sense of departure
1141 from homeostasis (i.e. the minimization of variational free energy) contingent upon remem-
1142 bered action-perception associations can always be expressed as a mixture of “epistemic”
1143 and “pragmatic” value. The pragmatic value is the expected outcome according to prior
1144 preferences, i.e. “good” or “bad” evaluations, while the epistemic value is the utility of the
1145 action for learning, i.e. reducing the potential for uncertainty or surprise in the future. This
1146 resolution of uncertainty through active sampling is at the heart of many active inference

1147 schemes and arises naturally in any model in which the agent expects to occupy the states
1148 it prefers.

1149 4.4 Reference frames and attention

1150 While defining expectations over percepts can be expected to be useful in some circum-
1151 stances, many aspects of realistic behavior require defining and acting on expectations
1152 defined over individual or small subsets of *components* of percepts. The memories X_{MD}
1153 and X_{MPA} together provide the data needed to allow individual component - action as-
1154 sociations to be computed; the memory X_E similarly provides the data needed to allow
1155 individual component valuations to be computed. Let X_C and X_{EC} be memories that store
1156 conditional probability distributions and evaluations, respectively, of individual components
1157 of percepts. To define X_C , note that the $x_R - g$ and $g - x_P$ associations stored in X_{MD} and
1158 X_{MPA} respectively allow each action g to be viewed as a relation $\{(x_R, x_P)\}$ implemented
1159 by PA . Expressing these percepts as vectors $x_R(t) = \sum_i \alpha_i(t)\xi_i$ and $x_P(t) = \sum_i \beta_i(t)\xi_i$,
1160 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
1161 can, in other words, be viewed as increasing or decreasing the amplitude of each percep-
1162 tual component ξ_i from one percept to the next. As it is natural to view amplitudes as
1163 probabilities of occurrence as discussed above, each g can be viewed as increasing or de-
1164 creasing the probability of each perceptual component ξ_i from one percept (i.e. value of
1165 t) to the next. The memory X_C can, therefore, be viewed as storing t -indexed conditional
1166 probabilities $\text{Prob}_t(\xi_i|g, \text{Prob}_{t-1}(\xi_i))$ of perceptual components given actions. To update
1167 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
1168 C as a map $X_{MD} \times X_{MPA} \times X_C \rightarrow X_C$. Subject to the constraint that all probabilities
1169 remain normalized, this map can in principle implement any arbitrary updating function.

1170 The memory X_{EC} containing component valuations may be constructed from X_E in a sim-
1171 ilar fashion, by defining punctual, forgetful kernels EC_{good} and EC_{bad} that map $X_E \rightarrow$
1172 X_{EC} . The kernels EC_{good} and EC_{bad} assign, respectively, “good” valuations to components
1173 strongly represented in “good” percepts and “bad” valuations to components strongly rep-
1174 resented in “bad” percepts. A suitable function for each would assign to each component
1175 ξ_i the average valuation of percepts x_P in which the coefficient α_i of ξ_i is greater than
1176 some specified threshold. With additional memory, this mechanism can be extended to
1177 assign values to (finite ranges of) amplitude values of components. Note that component
1178 valuations constructed in this way are in an important sense context-free; representing com-
1179 ponent valuations conditioned on the valuations of other components requires both more
1180 memory and more complex kernels.

1181 The memory components X_C and X_{EC} provide the “background knowledge” required for
1182 component-directed as opposed to entire-percept directed actions. What remains to be
1183 constructed is a process of selecting a component on which to act, and a second component
1184 with respect to which the action is taken. Consonant with current usage in physics (e.g.
1185 Bartlett, Rudolph and Spekkens, 2007), we refer to this second, context-setting component
1186 as a reference frame for the action. Specifying a reference frame is specifying what does

1187 *not* change when an action is taken; hence reference frames provide the basis for specifying
1188 what does change. Reference frames provide, in other words, the necessary stasis with
1189 respect to which change is perceptible. Measurement devices such as meter sticks provide
1190 the canonical example: a measurement made with a meter stick is only meaningful if one
1191 assumes that the actions involved in making the measurement do not change the length
1192 of a meter stick. More broadly, any context in which observations are made, whether a
1193 particular laboratory set-up or an everyday scene, is meaningful as a context only if it
1194 itself does change as a result of making the observation. A reference frame is, therefore, a
1195 *stipulated* solution to the frame problem, the problem of specifying what does not change
1196 as a result of an action (McCarthy and Hayes, 1969; reviewed by Fields, 2013b). Such
1197 stipulations are inherently fragile and defeasible: a context that does observably change,
1198 like a “meter stick” with an observably context-dependent length, ceases to be a reference
1199 frame as soon as its variation is detected. Stipulated reference frames are, nonetheless,
1200 *useful* solutions to the frame problem to the extent that they enable successful behavior in
1201 the niche of the agent employing them. Absent a level of control over the environment that
1202 ITP forbids, they are the only kinds of reference frames available.

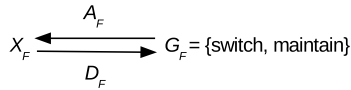
1203 While the frame problem has a long history in AI, its impact on cognitive science more
1204 generally has been primarily philosophical (see, e.g. the contributions to Pylyshyn (1987)
1205 and Ford and Pylyshyn (1996)). The question of how human perceivers identify *contexts*
1206 as opposed to objects or events and how they detect changes in context have received little
1207 direct investigation. The current model predicts that contexts are defined constructively
1208 by the activation of discrete reference frames that impose expectations of constancy and
1209 limit attention to features expected to remain constant. Experimental demonstrations of
1210 change-blindness (reviewed by Simons and Ambinder, 2005) show that such limitations of
1211 attention exist. Virtual reality methods provide opportunities to experimentally manipulate
1212 context identification, and hence to probe the specific reference frames employed to identify
1213 contexts, in ways that remain largely unexplored.

1214 For complex organisms, the most important reference frame is arguably the *experienced self*,
1215 generally including one or more distinguishable components of the *body*. This experienced
1216 self reference frame comprises a collection of components of experience that do not change
1217 during some, most or even all actions. The experienced self as a reference frame appears to
1218 be innate in humans (e.g. Rochat, 2012) and may be innate in higher animals generally. It is
1219 with respect to the experienced self as a reference frame that infants learn their capabilities
1220 for actions as bodily motions and for social interactions as communications with others (e.g.
1221 von Hofsten, 2007). Actions of or on the body, e.g. moving a limb, require that other parts
1222 of the experienced self, e.g. the mass and shape of the limb and its point of connection
1223 to the rest of the body, remain fixed to serve as the reference frame for the action. As
1224 the body grows and develops, its representation must be updated to compensate for these
1225 changes if its function as a reference frame is to be preserved. The experienced self reference
1226 frame is readily extensible to tools, vehicles, and fully-virtual avatars in telepresence and
1227 virtual-reality applications, and is readily manipulated in the laboratory. Disruptions of the
1228 experienced self as a reference frame present as pathologies ranging from schizophrenia to

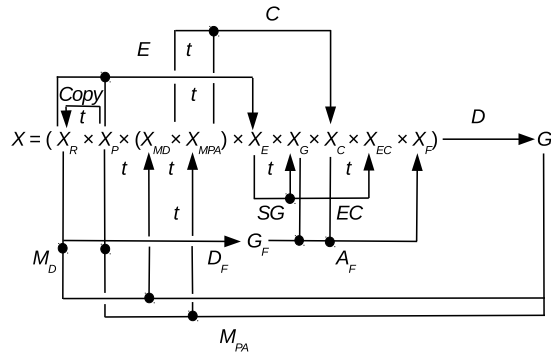
1229 anosognosia. These latter provide a clinical window into the human implementation of the
 1230 bodily and emotive self as a fusion of interoceptive and perceptual inputs (e.g. Craig, 2010;
 1231 Seth, 2013) and of the cognitive self as a fusion of memory-access and executive functions
 1232 that develops gradually from infancy to early adulthood (e.g. Simons, Henson, Gilbert and
 1233 Fletcher, 2008; Metzinger, 2011; Hohwy, 2016).

1234 Selecting a particular component of a percept on which to act and another component
 1235 or components, such as the experienced self or the experienced self in some perceived
 1236 surroundings, to serve as a fixed context for an action is an act of *attention*. The selected
 1237 components must, moreover, remain subjects of attention throughout the action. Any
 1238 agent capable of attending to some component of an ongoing scene must also, however,
 1239 be capable of switching attention to a different component if something unexpected and
 1240 important happens. Attention requires, therefore, not just a decision about what to attend
 1241 to, but also a decision about whether to maintain or switch attentional focus. To meet these
 1242 requirements, we introduce an “attentional workspace” X_F , a memory that contains a goal-
 1243 dependent focus of attention ξ_i , a focus-dependent reference frame ξ_j and a time counter
 1244 t_F that measures the duration of an attentional episode. We also define an attentional
 1245 action space G_F containing two actions, ‘switch’ and ‘maintain’ that alter or preserve the
 1246 attentional focus, respectively, and a forgetful punctual kernel $D_F : X_P \times X_R \times X_E \times X_G \rightarrow$
 1247 G_F that selects $g_F = \text{‘switch’}$ at t if the valuation of $x_P(t)$ differs from that of $x_R(t)$ by
 1248 some specified threshold and selects $g_F = \text{‘maintain’}$ otherwise. These elements of G_F
 1249 correspond to actions A_F on the workspace X_F , as shown in Fig. 10a. The action A_{Fm}
 1250 selected by $g_F = \text{‘maintain’}$ only increments t_F . The action A_{Fs} selected by $g_F = \text{‘switch’}$
 1251 selects a new focus of attention ξ_k , a new reference frame ξ_l and resets t_F to zero. We
 1252 represent this action as a forgetful punctual kernel $A_{Fs} : X_P \times X_G \times X_C \times X_{CE} \rightarrow X_F$.
 1253 How this attention-switching kernel is defined has a potentially large impact on the behavior
 1254 of the RCA whose attentional workspace X_F it affects. A rational A_{Fs} could be expected
 1255 to select a component ξ_i on which to focus that had a relatively large amplitude α_i in
 1256 both the current percept x_P and a high-value goal and a reference frame ξ_j , also with a
 1257 relatively large amplitude in both x_P and the goal, that was affected in the past primarily
 1258 by actions that did not affect ξ_i . While the valuation of the attentional focus ξ_i may be
 1259 “bad,” a rational A_{Fs} would select a reference frame ξ_j with a “good” or at least not “bad”
 1260 valuation, as this amplitude of this component is meant to be kept fixed in subsequent
 1261 interactions with W . A rational D kernel acting on the workspace X_F would then choose
 1262 actions g that, in the past as recorded in X_C , moved the amplitude of x_i in the direction of
 1263 its value in the chosen goal state while keeping the amplitude of x_j fixed. As X_C , X_{EC} and
 1264 X_F are updated one cycle behind X_{MD} , X_{MPA} , X_E and X_G and hence two cycles behind
 1265 X_P , the kernel D must always work with expectation and valuation information that is
 1266 slightly out-of-date.

a)



b)



1267 *Fig. 10:* a) Kernels that maintain or switch attentional focus. b) Additions to
 1268 *Fig. 9* required to support attention. Connections to W are again elided for
 1269 clarity.

1270 The structure of and operations within the experiential space X required for an atten-
 1271 tional system are summarized in *Fig. 10b*. Selecting a new component for attention and
 1272 maintaining attention on a previously-selected component are competitive processes in this
 1273 architecture, as they are in humans (reviewed by Vossel, Geng and Fink, 2014). When
 1274 top-down goals and expectations dominate and hence the dorsal attention system controls
 1275 perceptual processing, the salience of goal-irrelevant stimuli is reduced; a switch to vigilance
 1276 and hence ventral attentional control, in contrast, reduces the salience of goal-relevant stim-
 1277 uli. Top-down, dorsal attentional dominance facilitates exploration and information gather-
 1278 ing, while bottom-up, ventral attentional dominance facilitates threat avoidance. This
 1279 attention switch can be incorporated into predictive coding and active inference models
 1280 using the concept of “precision” for both expectations and percepts; high-precision expect-
 1281 ations dominate low-precision percepts and vice-versa (Friston, 2010; 2013). Precision is
 1282 effectively a measure of reliability based on prior experiences and is hence a second-order
 1283 expectation that must be learned by refining an *a priori* bias as discussed above. Predic-
 1284 tive coding networks modulated by estimated precision have been shown to describe the
 1285 cellular-scale connection architecture of cortical minicolumns (Bastos et al., 2012) as well
 1286 as the modular connection architectures of motor (Shipp, Adams and Friston, 2013) and
 1287 visual (Kanai, Komura, Shipp and Friston, 2015) processing (see also Adams, Friston and
 1288 Bastos (2015) for an overview of these results). As noted earlier, the smoothness of stored
 1289 probability distributions provides a natural estimate of the number of experiences that have
 1290 contributed to them and hence their reliability. A rational switching function can be ex-
 1291 pected to favor high-reliability expectations and disfavor low-reliability expectations, and
 1292 hence to implement a precision-based modulation of attention.

1293 Extending the system shown in *Fig. 10b* to multiple focus and/or reference components
 1294 costs memory and processing complexity, but does not change the architecture. It is inter-

1295 esting to note that within this architecture, all change is implicitly attributed by the agent
1296 to the action taken; from the agent’s intrinsic perspective, its actions change the state of its
1297 attentional focus with respect to its reference frame. For the system to behave effectively,
1298 the world W must be such that this attribution of observed changes to executed actions is
1299 satisficing in W . The world must not, in other words, surprise the agent so often that the
1300 agent’s sense that actions have predictable consequences becomes impossible to maintain.
1301 The world must not, in other words, exhibit either overall randomness or overall stasis as
1302 noted earlier.

1303 It is worth re-emphasizing, moreover, that in the CA framework X is a space of *experi-*
1304 *ences*. Hence the RCA depicted in Fig. 10b is regarded as *experiencing* each state of its
1305 highly-structured space X , including all those components on which its attention is *not*
1306 focussed (the formalism leaves open the question of whether these components themselves
1307 have unexperienced internal structure). It may, however, be “unconscious” of unattended
1308 components in the sense in which this term is used in theories that associate consciousness
1309 with relative amplification or attention (e.g. Baars, Franklin and Ramsoy, 2013; Dehaene,
1310 Charles, King and Marti, 2014; Graziano, 2014). In general, how an RCA acts depends
1311 on its attentional focus. Reporting what it is experiencing, e.g. to an investigator in a
1312 laboratory or even to itself via a modality such as inner speech, is a specific kind of ac-
1313 tion that requires a specific attentional focus. Whether the attentional focus required to
1314 support a given form of reporting is achieved in any particular case or is even achievable
1315 by a particular RCA is a matter of architecture, i.e. of how the memory-construction and
1316 attentional-control kernels are defined. Agents that never report particular kinds of experi-
1317 ences, or that never report experiences using a given modality such as inner speech (Heavey
1318 and Hurlburt, 2008), are not only possible but to be expected within the CA framework.
1319 Indeed the CA framework predicts that *agents are typically aware of more than they can*
1320 *report awareness of* to an external observer or even to themselves. Agents are, in other
1321 words, typically *under-equipped with attentional resources*, and hence unable to access some
1322 or even much of their experience for behavioral reporting via any particular modality. Being
1323 under-equipped for reporting experiences *post hoc* is unsurprising on evolutionary grounds;
1324 indeed why human beings should engage in so much *post hoc* self-reporting via modalities
1325 such as inner speech remains a mystery (Fields, 2002). As reportability by some observable
1326 behavior remains the “gold standard” in assessments of awareness (e.g. Dehaene, Charles,
1327 King and Marti, 2014), this strong and counter-intuitive prediction of the CA framework
1328 can at present only be tested indirectly, e.g. using phenomena such as blindsight (re-
1329 viewed by Overgaard, 2011). It raises the methodological question of whether “reporting”
1330 of experiences by imaging methods such a fMRI, as employed by Boly, Sanders, Mashour
1331 and Laureys (2013), for example, with otherwise-unresponsive coma patients, should be
1332 regarded as evidence of awareness across the board.

1333 4.5 Remembering and planning action sequences

1334 The attentional workspace X_F defined above does not explicitly represent the action taken
1335 at each t and so cannot support either memory for “cases” of successful action or plan-
1336 ning. The most recently executed g is, however, available within X_{MD} . A fixed-capacity
1337 case memory can be regarded as a subjective probability distribution over possible cases,
1338 where each case is a vector of fixed length l_{case} , the components of which are quadruples
1339 $(\alpha_i \xi_i, \beta_j \xi_j, t_F, g(t_F))$ with the percept components ξ_i, ξ_j and the amplitude β_j fixed. A case
1340 defined in this way provides a representation of how the amplitude α_i of the attentional
1341 focus ξ_i varies relative to the fixed amplitude β_j of the reference frame ξ_j when subjected
1342 to the sequence $g(t_F = 0) \dots g(t_F = l_{case})$ of actions. This definition formulates in lan-
1343 guage compliant with ITP the concept of a case employed in the case-based reasoning and
1344 planning literature (Riesbeck and Schank, 1989; Kolodner, 1993). It is also similar in both
1345 role and scope to the concept of an “event file” introduced by Hommel (2004) to repre-
1346 sent the temporal binding of perceptions with context-appropriate actions. Cases or event
1347 files are effectively “snapshots” of active inference that show how a particular perceptual
1348 input is processed given the attentional context in which it is received and the particular
1349 expectations that it activates.

1350 As an example, consider a sequence of actions involved in reaching for and grasping a coffee
1351 cup. The immediate goal of the sequence is to grasp the coffee cup; we will ignore the
1352 question of different grasps being needed for different subsequent actions. The target of
1353 the sequence is a *particular* coffee cup that is visually identifiable by particular perceived
1354 features, e.g. location, size, shape and color. The cup’s perceived size, shape and color do
1355 not change as a result of the motion; hence their values can serve as the reference frame
1356 that determines the cup’s identity. As the goal of the action sequence is to change the
1357 perceived location of the coffee cup, its location cannot be included in the reference frame;
1358 if it was, the cup would lose its identity when it was moved. The attentional workspace
1359 X_F , therefore, contains the variable perceived values of the positions of the cup and of the
1360 reaching hand as foci and the fixed perceived values of the size, shape and color of the cup
1361 as the reference frame. The recorded case contains, effectively, a sequence of “snapshots”
1362 of the contents of X_F : a time sequence of cup and hand position values, together with the
1363 actions that produced them, relative to these fixed reference values. A memory M_{case} for
1364 such cases can be constructed using the counter-incrementing methods used to construct
1365 X_{MD} and X_{MPA} above. As action sequences that are worth recording are typically those
1366 that either satisfied goals or led to trouble, it is useful to construct each record in M_{case} as
1367 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
1368 $x_P(t_F = 0)$ and $x_P(t_F = l_{case})$ are the full percepts at the beginning and the end of $case(t_F)$
1369 respectively, and $E((x_P(t_F = 0)))$ and $E((x_P(t_F = l_{case})))$ are their evaluations as recorded
1370 in X_E . This representation allows M_{case} to be searched – i.e. kernels acting on M_{case} to
1371 depend upon – either the initial state and its evaluation or the final state and its evaluation.
1372 Case memories constructed in this way are clearly combinatorially explosive; hence case-
1373 based planning in systems with limited memory is necessarily heuristic, not exhaustive, a
1374 condition widely recognized in the case-based planning literature.

1375 It is natural to interpret a set of one or more fixed components of experience, with respect
1376 to which one or more other components of experience change when one or more sequences
1377 of actions is executed as defining an effective or apparent *object*. Objects defined in this way
1378 are collections of expectations, based on accumulated experience, about the co-occurrence
1379 and co-variation under actions of particular values of particular experiential degrees of
1380 freedom. Objects in this sense are effectively *categories* defined by fixed (i.e. reference) and
1381 variable features together with sets of expected behaviors, i.e. changes in the amplitudes of
1382 the variable features relative to the fixed features in response to actions. Hence such objects
1383 are more properly considered to be object *types* as opposed to *de re* individuals. While
1384 an agent may *assume*, as a useful heuristic, that an object category has only one member
1385 and act on the basis of this assumption, consistency with ITP requires that nothing in
1386 the agent’s experience can be sufficient to demonstrate that this is the case. Hence object
1387 identity over time is ambiguous in principle in the ITP/CA framework. Objects defined in
1388 this way play the role of “icons” on the ITP interface. As the number of recorded cases
1389 involving actions that change the state of some object increase, its “icon” gains predictable
1390 functionality and hence utility as a locus of behavior.

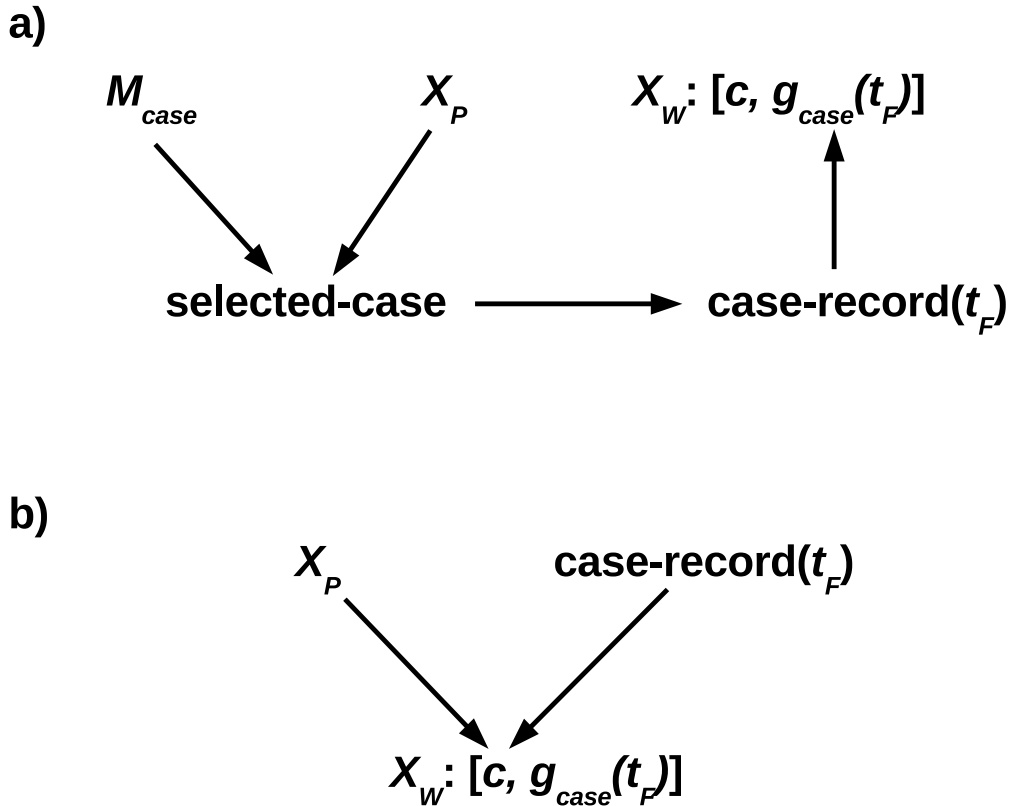
1391 The present framework leaves open the question of whether any “object”-specifying ref-
1392 erence frames are innate. It predicts, however, that any such reference frames, whether
1393 innately specified or constructed from experience, will have low dimensionality compared
1394 to the perceptual experiences that they help to interpret. Dramatic evidence for low di-
1395 mensionality is provided by studies of two of the earliest-developing and ecologically most
1396 crucial reference frames for humans, those that identify animacy and agency (reviewed by
1397 Scholl and Tremoulet, 2000; Scholl and Gao, 2013; Fields, 2014). Indeed Gao, McCarthy
1398 and Scholl (2010) have shown that a simple oriented “V” shape not only satisfies the typical
1399 human visual criterion for agency detection, but distracts attention sufficiently to disrupt
1400 performance in an object-tracking task. Human face-recognition criteria are similarly rudi-
1401 mentary. Additional evidence for low reference-frame dimensionality is provided by the
1402 kinds of categorization conflicts studied in the quantum cognition literature (reviewed e.g.
1403 by Pothos and Busemeyer, 2013; Bruza, Kitto, Ramm and Sitbon, 2015), for example the
1404 “Linda” problem. Here the “natural” reference frames, i.e. concepts or coherent sets of
1405 expectations, do not exhibit classical compositionality; combining reference frames to repro-
1406 duce the judgements made by subjects requires the use of complex “quantum” probability
1407 amplitudes. Complex probabilities can, however, be represented by classical probabilities
1408 in higher-dimensional spaces (e.g. Fuchs and Schack, 2013; see also Fields, 2016 for a less
1409 formal discussion), consistent with attentional selection of a low-dimensional subspace to
1410 serve as a reference frame. If “object”-specifying reference frames in fact encode fitness
1411 information as ITP requires, one would expect a general inverse correlation between fitness
1412 consequences and reference frame dimensionality. While both the global and local struc-
1413 ture of the typical human category hierarchy have been investigated (reviewed by Martin,
1414 2007; Keifer and Pulvermüller, 2012), neither the minimal functional content (i.e. dimen-
1415 sionality) nor the fitness-dimensionality correlation of typical categories have been broadly
1416 investigated.

1417 The components of the experienced self reference frame, taken together, constitute an
1418 iconic object – the experienced self as a persistent embodied actor – in the above sense.
1419 The features of the experienced self as persistent embodied actor that are employed as
1420 fixed reference features with respect to which other features of the experienced self are
1421 allowed to vary change only slowly and asynchronously as a function of time; it is this slow
1422 and asynchronous change in reference features that allow the approximation of a persistent
1423 experienced self (but see Klein, 2014 for a discussion of the sense of a persistent experienced
1424 self in the presence of conflicting perceptual evidence). The conditions under which non-
1425 self objects are represented as persistent over extended time, in particular across extended
1426 periods of non-observation, have been subjected to surprisingly little direct experimental
1427 investigation and are not well understood (e.g. Scholl, 2007; Fields, 2012). Both the
1428 extensibility of the experienced self reference frame to incorporate otherwise non-self objects
1429 discussed earlier and the sheer variety of pathologies of the experienced self, including
1430 depersonalization syndromes (e.g. Debruyne, Portzky, Van den Eynde and Audenaert,
1431 2009), suggest that the experienced self - non-self distinction is not constant for individual
1432 human subjects and highly variable between subjects. This question cannot, unfortunately,
1433 yet be addressed productively in non-human subjects.

1434 With this concept of an iconic object, the functional difference between a case memory
1435 M_{case} and the event memories X_{MD} and X_{MPA} becomes clear: M_{case} records sequences of
1436 *partial* events in which, in each sequence, only the response to actions of the attentional
1437 focus ξ_i and the lack of response to actions of the reference ξ_j are made explicit. Each case
1438 in M_{case} can, therefore, be thought of as imposing an implicit, goal-dependent criterion of
1439 *relevance* on the actions it records.

1440 Recording object-directed action sequences is useful to an agent because it enables previously-
1441 successful sequences to be repeated and previously-unsuccessful sequences to be avoided.
1442 Selecting a previously-recorded case from memory for execution under some similar cir-
1443 cumstances is the simplest form of planning. Executing the action sequence recorded in a
1444 remembered case requires, however, shortcutting the usual decision process D . Within the
1445 architecture shown in Fig. 10, the simplest way to accomplish this is to associate a working
1446 memory X_W with the attentional focus X_F , and to include in X_W a control bit c on which
1447 D depends. If $c = 0$, D is independent of the contents of X_W and acts as in Fig. 9. If $c = 1$,
1448 D selects the action g represented in X_W . Populating X_W requires two embedded agents,
1449 as shown in Fig. 11. The first agent (Fig. 11a) selects a recorded case based on the current
1450 percept, and sequentially copies the actions specified by that case into X_W . The “world” of
1451 this agent consists of X_P , M_{case} and X_W ; its “perception” kernel selects the case from M_{case}
1452 for which the initial state is closest to the current percept x_P , its “decision” kernel selects
1453 records from this case in sequence and its “action” kernel writes the action $g(t_F)$ specified
1454 by the selected case into X_W . The process executed by this agent requires a time step,
1455 i.e. one increment of t . The second agent (Fig. 11b) has a switching function analogous
1456 to the attention-switching dyad in Fig. 10a: it compares the current percept $x_P(t)$ to the
1457 currently-selected case record, setting $c = 1$ when the case is initially selected and setting
1458 $c = 0$ if the distance between the states of either the object or reference components of $x_P(t)$

1459 and their states as specified by the currently-selected case record exceeds some threshold.
 1460 Setting $c = 0$ in response to such an expectation violation during case execution restores
 1461 D to its usual function. Maintaining temporal synchrony requires that the overall counter
 1462 t advances only when D executes as discussed above; this requirement can be met if D is
 1463 regarded as acting instantaneously when $c = 1$ and the action g to be selected is specified
 1464 by X_W , i.e. when action is performed “automatically.” In this case interrupting execution
 1465 of a case must be regarded as requiring one time step, after which no action is selected.



1466 *Fig. 11:* a) Selection of a case and case-record for execution based on the current
 1467 percept. This action does not enable case execution. b) Enabling or disabling
 1468 case execution by setting or resetting the control bit c based on a comparison
 1469 of current and expected percepts during case execution.

1470 The processes illustrated in Fig. 11 only execute a previous case verbatim. Interrupting
 1471 execution of a case initiates a search for a new case that is a better fit to the current per-
 1472 cept $x_P(t)$. A more intelligent case-based planner can be constructed by incorporating an
 1473 additional agent capable of modifying the currently-selected case record based on $x_P(t)$ and
 1474 information about previous component responses stored in X_C . Such modification creates

1475 a new case, which is then recorded in M_{case} . A second natural extension would incorporate
1476 a “meta” agent capable of comparing multiple cases to identify shared perception-action
1477 dependencies. A case comparator of this kind is the minimal structure needed to recognize
1478 relationships between events occurring in different orders or with different numbers of in-
1479 tervening events; hence it is the minimal structure needed to implement a “temporal map”
1480 as described by Balsam and Gallistel (2009).

1481 5 Conclusion

1482 We have shown three things in this paper. First, the CA formalism introduced by Hoffman
1483 and Prakash (2014) is both powerful and non-trivial. Even “agents” comprising only a
1484 handful of bits exhibit surprisingly complex behavior. A three-bit agent can implement a
1485 Toffoli gate, so networks of three-bit agents can compute any computable function, and
1486 can even do so reversibly. More intriguing are the hints that networks of simple agents
1487 exhibit dynamical symmetries that also characterize geometry. This result comports well
1488 with current efforts by physicists to derive the familiar geometry of spacetime from the
1489 symmetries of information exchange between simple processing units (e.g. Tegmark, 2015).
1490 We are currently working toward a full description of spacetime constructed entirely within
1491 the CA framework.

1492 We have, second, shown that concept of “fitness” as connectivity emerges naturally when
1493 networks of interacting RCAs are considered. This fitness concept accords well with estab-
1494 lished concepts of centrality developed in the theory of social and other complex networks.
1495 By expressing fitness with the CA framework, we free ITP from any need to rely on an
1496 externally-stipulated fitness function. Computational experiments to characterize the con-
1497 ditions in which preferential attachment and hence high-connectivity individuals emerge in
1498 networks of interacting RCAs are being designed.

1499 Our third result is that networks of RCAs can, at least in principle, implement sophisticated
1500 cognitive processes including attention, categorization and planning. This result fleshes out
1501 the central concepts of ITP: that experience is an *interface* onto an ontologically-ambiguous
1502 world, and that “objects” and “causal relations” are patterns of positive and negative cor-
1503 relations between experiences. It highlights the critical role played by aspects of experience
1504 that do not change, and hence serve as “context” or, more formally, reference frames rel-
1505 ative to which aspects of experience that do change can be classified and analyzed. Here
1506 again, our result comports well with recent work in physics, where with the rise of quantum
1507 information theory, the roles of reference frames in defining what can and cannot be known
1508 or communicated about a physical situation have taken on new prominence (e.g. Bartlett,
1509 Rudolph and Spekkens, 2007). A substantial program of simulation development and test-
1510 ing is clearly required to evaluate, in structured and eventually in open environments, the
1511 formal models of memory, attention, categorization and planning developed here. The level
1512 of complexity at which such models can feasibly be implemented remains unclear. We hope,
1513 however, to be able to fully characterize the reference frames required to support relatively

1514 simple behaviors in relatively simple environments, and to use this information to formulate
1515 predictions testable in more complex systems.

1516 The CA framework is, as we have emphasized, a minimal formal framework for under-
1517 standing cognition and agency. While debates about the structure and content of memory
1518 - and implicitly, experience - have dominated cognitive science for decades (e.g. Gibson,
1519 1979; Fodor and Pylyshyn, 1988; Anderson, 2003), these debates have generally been con-
1520 ducted either informally or in the context of complex, conceptually open-ended modeling
1521 paradigms. Our results, together with those of Friston and colleagues using the predictive
1522 coding and adaptive inference framework, show that cognition and agency can be addressed
1523 in conceptually very simple terms. The primary task of an organism in an environment
1524 is to regulate its interactions with the environment, by behaving appropriately, in order
1525 to maintain an environmental state conducive its own homeostasis. As Conant and Ashby
1526 (1970) showed and Friston (2010; 2013) has significantly elaborated, effective regulation
1527 of the environment requires a statistically well-founded model of the environment. Consis-
1528 tency with ITP requires that such models treat the environment as open, in which case they
1529 can be at best satisficing. The results obtained here, together with those of Friston (2013)
1530 and Friston, Levin, Sengupta and Pezzulo (2015), offer an outline of how such models may
1531 be constructed in a way that is consistent with ITP, but many details remain to be worked
1532 out. A thorough treatment of both evolutionary and developmental processes from both
1533 extrinsic and intrinsic perspectives is needed to understand the kinds of worlds W in which
1534 complex networks of interdependent RCAs can be expected to appear.

1535 We have largely deferred the question of motivation. As mentioned in §4.3 above, ratio-
1536 nal agents exhibit curiosity and hence explore their environments to discover sources of
1537 “good” experiences, which in a typical W may lie very near sources of “bad” experiences.
1538 As Gottlieb, Oudeyer, Lopes and Baranes (2013) emphasize, however, rational agents do
1539 not exhibit unlimited curiosity, as this can lead to expending all available resources at-
1540 tempting to solve unsolvable problems or learn unlearnable information. Understanding
1541 and modeling motivation requires not only a formal characterization of resources and their
1542 use, but also a formal model of reward, its representation, and its roles in both extrinsic
1543 and intrinsic motivation. The distinction between the “pragmatic” and “epistemic” values
1544 of information (Friston et al., 2015) is useful here; the current framework models the effects
1545 of this distinction in terms of attention switching, but not its origin. Both developmental
1546 robotics (e.g. Cangelosi and Schlesinger, 2015) and the neuroscience of the reward system
1547 (e.g. Berridge, and Kringelbach, 2013) provide empirical avenues to pursue in this regard.

1548 We have also, and more importantly from an architectural perspective, deferred the task
1549 of constructing a full theory of RCA networks and RCA combinations. Developing such
1550 a theory will require addressing such questions as whether RCA networks can in general
1551 be considered locally hierarchical, whether the action spaces G of complex RCAs require
1552 structures, for example to represent fully automatized action patterns, analogous to the
1553 structures in X described here, and how to explicitly define D kernels in complex RCAs. It
1554 will also require understanding how the time counters (i.e. t parameters) of complex RCAs
1555 relate to those of their component RCAs, a question that has been elided here by assuming

1556 that all processes “inside” X are synchronous. Answering such questions may well depend
1557 on resolving at least some of the issues having to do with fitness and motivation mentioned
1558 above. We expect, however, that their answers will shed light on such questions as whether
1559 complex RCAs can in some cases be regarded as unaware of the experiences - e.g. the
1560 percepts or memories - of their component RCAs and how the actions of complex RCAs
1561 depend, or not, on the actions of their component RCAs.

1562 As CAs and hence RCAs are intended, from the outset, to represent *conscious* agents, it
1563 is natural to ask what the behavior of networks of RCAs can tell us about consciousness.
1564 Here two results stand out. The first is that an agent cannot, without violating ITP,
1565 distinguish the world outside of her experience from another conscious agent. While this
1566 follows from the ontological principle of conscious realism of Hoffman and Prakash (2014),
1567 it equally follows from the impossibility, within ITP, of determining that the “world” has
1568 non-Markovian dynamics. The second is that agents can be expected to be aware of more
1569 than they can report. This seems paradoxical if awareness is equated with reportability,
1570 but makes sense when the attentional resources that would be required to enable reporting
1571 of all experiences are taken into account.

1572 While examining specific cases of successful and unsuccessful behavior in well-defined worlds
1573 requires addressing the issues of motivation and multi-agent combination highlighted above,
1574 two substantial conceptual issues stand out. The first is that the CA formalism, in contrast
1575 to either standard neural network approaches or purely-functional cognitive modelling ap-
1576 proaches, enforces by its structure a focus on what a constructed agent is being modelled
1577 as experiencing. The CA formalism itself requires that the decision kernel D acts on the
1578 space of experiences X ; hence whatever D acts on must be in X and therefore must be an
1579 experience. Constructing complex memory structures in X in order to make them available
1580 to D is, given this constraint, proposing the hypothesis that the contents of such struc-
1581 tures are experienced. Experienced by whom? Here the second issue becomes relevant.
1582 As discussed in §3.2, discussions of consciousness have often assumed, explicitly or more
1583 typically implicitly, that “low-level” experiences combine in some straightforward way into
1584 “higher-level” experiences. The phenomenal unity of ordinary, waking human experience
1585 is assumed by many to indicate that there is only one relevant “level” of experience, the
1586 level of the whole organism (or often, just its brain). With this assumption, proper com-
1587 ponents of the human neurocognitive system cannot themselves be experiencers; that this
1588 is the case is treated as axiomatic, for example, in Integrated Information Theory (Tononi
1589 and Koch, 2015; see Cerullo, 2015 for a critique of this assumption in the IIT context).
1590 If complex experiencers are networks of RCAs, however, this assumption cannot be cor-
1591 rect: all RCAs, even the simplest ones, experience *something*. If complex experiencers are
1592 networks of RCAs, there is also no reason to assume that “higher-level” experiences are in
1593 any straightforward sense combinations of “lower-level” ones. Unless RCA combinations are
1594 simple Cartesian products, high-level experiences will in general not be uniquely predictable
1595 from low-level experiences or vice-versa. If complex experiencers are only approximately
1596 hierarchical rich-club networks of RCAs, the assumption that experiences should in general
1597 be straightforwardly combinatoric is almost certainly wrong.

1598 That said, it is worth re-emphasizing that the CA framework is not, and is not intended to
1599 be, a theory of consciousness *per se*. The CA framework says nothing about the *nature* of
1600 experience. It says nothing about qualia; it simply assumes that qualia exist, that agents
1601 experience them, and that they can be tokened by elements of X . The CA framework is,
1602 instead, a formal framework for modelling conscious agents and their interactions that en-
1603 forces consistency with ITP. By itself, the CA framework is ontologically neutral, as is ITP.
1604 When equipped with the ontological assumption of conscious realism, the CA framework
1605 becomes at least *prima facie* consistent with ontological theories that take consciousness to
1606 be an irreducible primitive. The role of the CA framework in expressing the assumptions
1607 or results of such theories can be expected to depend on the details of their ontological
1608 assumptions. Whether the CA framework fully captures the ontological assumptions of
1609 existing theories that take consciousness to be fundamental, e.g. that of Faggin (2015),
1610 remains to be determined.

1611 In summary, the CA framework, and RCA networks in particular, provide both a highly-
1612 constrained formal technology for representing cognition and a way of thinking about cogni-
1613 tion that emphasizes experience and decisions based on experience. It directly implements
1614 the ontological neutrality regarding the external world that is required by ITP. As results
1615 from physics and other disciplines render naïve or even critical realism about perceived
1616 objects and causal relations increasingly hard to sustain, this ability to model experience
1617 and decision making with no supporting ontology will become increasingly critical for psy-
1618 chology and for the biosciences in general.

1619 Acknowledgements

1620 The authors thank Federico Faggin and Robert Prentner for discussions of the ideas in this
1621 paper and The Federico and Elvia Faggin Foundation for financial support. Thanks also
1622 to the reviewers for their constructive comments.

1623 References

- 1624 Adams, R. A., Friston, K. J. and Bastos, A. M. (2015). Active inference, predictive coding
1625 and cortical architecture. In M. F. Casanova and I. Opris (Eds) *Recent Advances in the*
1626 *Modular Organization of the Cortex*. Berlin: Springer (pp. 97-121).
- 1627 Adolphs, R. (2003). Cognitive neuroscience of human social behavior. *Nature Reviews*
1628 *Neuroscience* 4, 165-178.
- 1629 Adolphs, R. (2009). The social brain: Neural basis for social knowledge. *Annual Review of*
1630 *Psychology* 60, 693-716.
- 1631 Agrawal, H. (2002). Extreme self-organization in networks constructed from gene expression
1632 data. *Physical Review Letters* 89, 268702.

- 1633 Anderson, M. L. (2003). Embodied cognition: A field guide. *Artificial Intelligence* 149,
1634 91-130.
- 1635 Ashby, W. R. (1956). *Introduction to Cybernetics*. London: Chapman and Hall.
- 1636 Aspect, A., Dalibard, J. and Roger, G. (1982). Experimental test of Bell's inequalities using
1637 time-varying analyzers. *Physical Review Letters* 49, 1804-1807.
- 1638 Baars, B. J., Franklin, S. and Ramsoy, T. Z. (2013). Global workspace dynamics: Cortical
1639 "binding and propagation" enables conscious contents. *Frontiers in Psychology* 4, Article
1640 # 200.
- 1641 Baillargeon, R. (2008). Innate ideas revisited: For a principle of persistence in infants
1642 physical reasoning. *Perspectives on Psychological Science* 3, 2-13.
- 1643 Barabási, A.-L. (2009). Scale-free networks: A decade and beyond. *Science* 325, 412-413.
- 1644 Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *Science*
1645 286, 509-512.
- 1646 Barabási, A.-L. and Oltvai, Z. N. (2004). Network biology: Understanding the cell's func-
1647 tional organization. *Nature Reviews Genetics* 5, 101-114.
- 1648 Bartlett, S. D., Rudolph, T. and Spekkens, R. W. (2007). Reference frames, superselection
1649 rules, and quantum information. *Reviews of Modern Physics* 79, 555-609.
- 1650 Bassett, D. S. and Bullmore, E. (2006). Small world brain networks. *The Neuroscientist*
1651 12, 512-523.
- 1652 Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P. and Friston, K. J.
1653 (2012). Canonical microcircuits for predictive coding. *Neuron* 76, 695-711.
- 1654 Bennett, B. M., Hoffman, D. D. and Prakash, C. (1989). *Observer Mechanics: A Formal*
1655 *Theory of Perception*. Academic Press.
- 1656 Berridge, K. C. and Kringelbach, M. L. (2013). Neuroscience of affect: brain mechanisms
1657 of pleasure and displeasure. *Current Opinion in Neurobiology* 23, 294-303.
- 1658 Boly, M., Sanders, R. D., Mashour, G. A. and Laureys, S. (2013). Consciousness and
1659 responsiveness: Lessons from anaesthesia and the vegetative state. *Current Opinion in*
1660 *Anesthesiology* 26, 444-449.
- 1661 Börner, K., Sanyal, S. and Vespignani, A. (2007). Network science. *Annual Review of*
1662 *Information Science and Technology* 41, 537-607.
- 1663 Bruza, P. D., Kitto, K., Ramm, B. J. and Sitbon, L. (2015). A probabilistic framework
1664 for analysing the compositionality of conceptual combinations. *Journal of Mathematical*
1665 *Psychology* 67, 26-38.
- 1666 Burgess, P. W. and Wu, H-C. (2013). Rostral prefrontal cortex (Brodmann area 10):
1667 Metacognition in the brain. In D. T. Stuss and R. T. Knight (Eds.) *Principles of Frontal*
1668 *Lobe Function, 2nd Ed.* New York: Oxford University Press (pp. 524-534).

- 1669 Cangelosi, A. and Schlesinger, M. (2015). *Developmental Robotics: From Babies to Robots*.
1670 Cambridge, MA: MIT Press.
- 1671 Cannon, T. D. (2015). How schizophrenia develops: Cognitive and brain mechanisms
1672 underlying onset of psychosis. *Trends in Cognitive Science* 19, 744-756.
- 1673 Cerullo, M. A. (2015). The problem with Phi: A critique of Integrated Information Theory.
1674 *PLoS Computational Biology* 11, e1004286.
- 1675 Colizza, V., Flammini, A., Serrano, M. A. and Vespignani, A. (2006). Detecting rich-club
1676 ordering in complex networks. *Nature Physics* 2, 110-115.
- 1677 Conant, R. C. and Ashby, W. R. (1970). Every good regulator of a system must be a model
1678 of that system. *International Journal of Systems Science* 1, 89-97.
- 1679 Conway, J. and Kochen, S. (2006). The free will theorem. *Foundations of Physics* 36,
1680 1441-1473.
- 1681 Craig, A. D. (2010). The sentient self. *Brain Structure and Function* 214, 563-577.
- 1682 Cummins, R. (1977). Programs in the explanation of behavior. *Philosophy of Science* 44,
1683 269-287.
- 1684 Damasio, A. (1999). *The Feeling of What Happens: Body and Emotion in the Making of*
1685 *Consciousness*. Orlando, FL: Harcourt.
- 1686 Debruyne, H, Portzky, M., Van den Eynde, F. and Audenaert, K. (2009). Cotards syn-
1687 drome: A review. *Current Psychiatry Reports* 11, 197-202.
- 1688 Dehaene, S., Charles, L., King, J.-R. and Marti, S. (2014). Toward a computational theory
1689 of conscious processing. *Current Opinion in Neurobiology* 25, 76-84.
- 1690 Diestel, R. (2010). *Graph theory* (4th ed.). Berlin: Springer.
- 1691 Dunbar, R. I. M. (2003). The social brain: Mind, language and society in evolutionary
1692 perspective. *Annual Review of Anthropology* 32, 163-181.
- 1693 Dunbar, R. I. M., and Shultz, S. (2007). Evolution in the social brain. *Science* 317, 1344-
1694 1347.
- 1695 Eibenberger, S., Gerlich, S., Arndt, M., Mayor, M. and Txen, J. (2013). Matter-wave
1696 interference of particles selected from a molecular library with masses exceeding 10,000
1697 amu. *Physical Chemistry and Chemical Physics* 15, 14696-14700.
- 1698 Faggin, F. (2015). The nature of reality. *Atti e Memorie dell'Accademia Galileiana di*
1699 *Scienze, Lettere ed Arti, Volume CXXVII (2014-2015)*. Padova: Accademia Galileiana di
1700 *Scienze, Lettere ed Arti*.
- 1701 Fields, C. (2002). Why do we talk to ourselves? *Journal of Experimental & Theoretical*
1702 *Artificial Intelligence* 14, 255-272.
- 1703 Fields, C. (2012). The very same thing: extending the object token concept to incorporate
1704 causal constraints on individual identity. *Advances in Cognitive Psychology* 8, 234-247.

- 1705 Fields, C. (2013a). A whole box of Pandoras: Systems, boundaries and free will in quantum
1706 theory. *Journal of Experimental & Theoretical Artificial Intelligence* 25, 291-302.
- 1707 Fields, C. (2013b). How humans solve the frame problem. *Journal of Experimental &*
1708 *Theoretical Artificial Intelligence* 25, 441-456.
- 1709 Fields, C. (2014). Motion, identity and the bias toward agency. *Frontiers in Human*
1710 *Neuroscience* 8, Article # 597.
- 1711 Fields, C. (2016). Building the observer into the system: Toward a realistic description of
1712 human interaction with the world. *Systems* 4, Article # 32.
- 1713 Fodor, J. A. and Pylyshyn, Z. W. (1988). Connectionism and cognitive architecture: A
1714 critical analysis. *Cognition* 28, 3-71.
- 1715 Ford, K. M. and Pylyshyn, Z. W. (Eds.) (1996). *The Robot's Dilemma Revisited*. Norwood,
1716 NJ: Ablex.
- 1717 Friston, K. (2010). The free-energy principle: A unified brain theory? *Nature Reviews*
1718 *Neuroscience* 11, 127-138.
- 1719 Friston, K. (2013). Life as we know it. *Journal of the Royal Society: Interface* 10, 20130475.
- 1720 Friston, K., Thornton, C. and Clark, A. (2012). Free-energy minimization and the dark-
1721 room problem. *Frontiers in Psychology* 3, article # 130.
- 1722 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P, O'Doherty, J. and Pezzulo, G.
1723 (2016). Active inference and learning. *Neuroscience and Biobehavioral Reviews* 68, 862-879.
- 1724 Friston, K., Levin, M., Sengupta, B. and Pezzulo, G. (2015). Knowing ones place: A
1725 free-energy approach to pattern regulation. *Journal of the Royal Society: Interface* 12,
1726 20141383.
- 1727 Friston, K., Rigoli, F., Ognibene, D., Mathys, C., FitzGerald, T. and Pezzulo, G. (2015).
1728 Active inference and epistemic value. *Cognitive Neuroscience* 6, 187-214.
- 1729 Fuchs, C. A. and Schack, R. (2013). Quantum-Bayesian coherence. *Reviews of Modern*
1730 *Physics* 85, 1693-1715.
- 1731 Gao, T., McCarthy, G. and Scholl, B. J. (2010). The wolfpack effect: Perception of animacy
1732 irresistibly influences interactive behavior. *Psychological Science* 21, 1845-1853.
- 1733 Gibson, J. J. (1979). *The Ecological Approach to Visual Perception*. Boston: Houghton-
1734 Mifflin.
- 1735 Geisler, W.S. and Diehl, R.L. (2003). A Bayesian approach to the evolution of perceptual
1736 and cognitive systems. *Cognitive Science* 27, 379-402.
- 1737 Giustina, M., Versteegh, M. A. M., Wengerowsky, S. et al. (2015). A significant-loophole-
1738 free test of Bells theorem with entangled photons. *Physical Review Letters* 115, 250401.
- 1739 Goldberg, R. P. (1974). A survey of virtual machine research. *IEEE Computer* 7(6), 34-45.

- 1740 Gottlieb, J., Oudeyer, P.-Y., Lopes, L. and Baranes, A. (2013). Information-seeking, curios-
1741 ity, and attention: Computational and neural mechanisms. *Trends in Cognitive Sciences*
1742 17, 585-593.
- 1743 Graziano, M. S. A. (2014). Speculations of the evolution of awareness. *Journal of Cognitive*
1744 *Neuroscience* 26, 1300-1304.
- 1745 He, X., Feldman, J. and Singh, M. (2015). Structure from motion without projective con-
1746 sistency. *Journal of Vision* 15, 725.
- 1747 Heavey, C. L. and Hurlburt, R. T. (2008). The phenomena of inner experience. *Conscious-*
1748 *ness and Cognition* 17, 798-810.
- 1749 Hensen, B., Bernien, H., Drau, A. E. et al. (2015). Loophole-free Bell inequality violation
1750 using electron spins separated by 1.3 kilometres. *Nature* 526, 682-686.
- 1751 Hoffman, D. D. (2016). The interface theory of perception. *Current Directions in Psycho-*
1752 *logical Science* 25, 157-161.
- 1753 Hoffman, D. D. and Prakash, C. (2014) Objects of consciousness. *Frontiers in Psychology*
1754 5, Article # 577.
- 1755 Hoffman, D. D. and Singh, M. (2012). Computational evolutionary perception. *Perception*
1756 41, 1073-1091.
- 1757 Hoffman, D. D., Singh, M. and Prakash, C. (2015). The interface theory of perception.
1758 *Psychonomic Bulletin & Review* 22, 1480-1506.
- 1759 Hohwy, J. (2016). The self-evidencing brain. *Noûs* 50, 259-285.
- 1760 Hommel, B. (2004). Event files: Feature binding in and across perception and action.
1761 *Trends in Cognitive Sciences* 8, 494-500.
- 1762 Jacques, V., Wu, E., Grosshans, F., Treussart, F., Grangier, P., Aspect, A. and Roch,
1763 J.-F. (2007). Experimental realization of Wheeler's delayed-choice gedanken experiment.
1764 *Science* 315, 966-968.
- 1765 Jennings, D. and Leifer, M. (2015). No return to classical reality. *Contemporary Physics*
1766 in press (arxiv:1501.03202).
- 1767 Kanai, R., komura, Y., Shipp, S. and Friston, K. (2015). Cerebral hierarchies: Predictive
1768 processing, precision and the pulvinar. *Philosophical Transactions of the Royal Society B*
1769 370, 20140169.
- 1770 Keifer, M. and Pulvermüller, F. (2012). Conceptual representations in mind and brain:
1771 Theoretical developments, current evidence and future directions. *Cortex* 7, 805-825.
- 1772 Kitto, K. (2014). A contextualised general systems theory. *Systems* 2, 541-565.
- 1773 Klein, S. B. (2014). Sameness and the self: Philosophical and psychological considerations.
1774 *Frontiers in Psychology* 5, Article # 29.
- 1775 Kolodner, J. (1993). *Case-Based Reasoning*. San Mateo, CA: Morgan Kaufmann.

- 1776 Koenderink, J. J. (2014). The all seeing eye? *Perception* 43, 1-6.
- 1777 Koenderink, J. J., van Doorn, A. J. and Todd, J. T. (2009). Wide distribution of external
1778 local sign in the normal population. *Psychological Research* 73, 14-22.
- 1779 Landauer, R. (1961). Irreversibility and heat generation in the computing process. *IBM J.*
1780 *Research Development* 5, 183-195.
- 1781 Landauer, R. (1999). Information is a physical entity. *Physica A* 263, 63-67.
- 1782 Landsman, N. P. (2007). Between classical and quantum. In: J. Butterfield and J. Earman
1783 (Eds) *Handbook of the Philosophy of Science: Philosophy of Physics*. Amsterdam: Elsevier.
1784 pp 417553.
- 1785 Lloyd, S. (2012). A Turing test for free will. *Philosophical Transactions of the Royal Society*
1786 *A* 370, 3597-3610.
- 1787 Maloney, L. T. and Zhang, H. (2010). Decision-theoretic models of visual perception and
1788 action. *Vision Research* 50, 2362-2374.
- 1789 Manning, A. G., Khakimov, R. I., Dall, R. G. and Truscott, A. G. (2015). Wheelers
1790 delayed-choice gedanken experiment with a single atom. *Nature Physics* 11, 539-542.
- 1791 Mark, J. T., Marion, B. B., and Hoffman, D. D. (2010). Natural selection and veridical
1792 perceptions. *Journal of Theoretical Biology* 266, 504-515.
- 1793 Marr, D. (1982). *Vision*. San Francisco: Freeman.
- 1794 Martin, A. (2007). The representation of object concepts in the brain. *Annual Review of*
1795 *Psychology* 58, 25-45.
- 1796 McCarthy, J. and Hayes, P. J. (1969). Some philosophical problems from the standpoint
1797 of artificial intelligence. In D. Michie and B. Meltzer (Eds.) *Machine intelligence, Vol. 4*.
1798 Edinburgh: Edinburgh University Press (pp. 463-502).
- 1799 Mermin, N. D. (1985). Is the moon there when nobody looks? Reality and the quantum
1800 theory. *Physics Today* 38(4), 38-47.
- 1801 Merton, R. K. (1968). The Matthew effect in science. *Science* 159, 56-63.
- 1802 Metzinger, T. (2011). The no-self alternative. In Gallagher, S (Ed.) *The Oxford Handbook*
1803 *of the Self*. Oxford: Oxford University Press (pp 287-305).
- 1804 Newman, M. E. J. (2001). The structure of scientific collaboration networks. *Proceedings*
1805 *of the National Academy of Sciences USA* 98, 404-409.
- 1806 Overgaard, M. (2011). Visual experience and blindsight: A methodological review. *Exper-*
1807 *imental Brain Research* 209, 473-479.
- 1808 Palmer, S. E. (1999). *Vision Science: Photons to Phenomenology*. Cambridge, MA: MIT
1809 Press.
- 1810 Parthasarathy, K. R. (2005). *Introduction to Probability and Measure*. Gurgaon, India:
1811 Hindustan Book Agency.

- 1812 Pattee, H. H. (2001). The physics of symbols: Bridging the epistemic cut. *Biosystems* 60,
1813 5-21.
- 1814 Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible*
1815 *Inference*. San Mateo CA: Morgan Kaufmann.
- 1816 Peil, K. (2015). Emotional sentience and the nature of phenomenal experience. *Progress*
1817 *in Biophysics and Molecular Biology* 119, 545-562.
- 1818 Pizlo, Z., Li, Y., Sawada, T. and Steinman, R.M. (2014). *Making a Machine that Sees Like*
1819 *Us*. New York: Oxford University Press.
- 1820 Polanyi, M. (1968). Lifes irreducible structure. *Science* 160, 1308-1312.
- 1821 Pont, S. C., Nefs, H. T., van doorn, A. J., Wijntjes, M. W. A., te Pas, S. F., de Ridder, H.
1822 and Koenderink, J. J. (2012). *Seeing and Perceiving* 25, 339-349.
- 1823 Popper, K. (1963). *Conjectures and Refutations: The Growth of Scientific Knowledge*.
1824 London: Routledge & Kegan Paul.
- 1825 Pothos, E. M. and Busemeyer, J. M. (2013). Can quantum probability provide a new
1826 direction for cognitive modeling? *Behavioral and Brain Sciences* 36, 255-327.
- 1827 Prakash, C. and Hoffman, D. D. (2016). Structure invention by conscious agents. In review.
- 1828 Prakash, C., Hoffman, D. D., Stephens, K. D., Singh, M. and Fields, C. (2016). Fitness
1829 beats truth in the evolution of perception. In review.
- 1830 Pylyshyn, Z. W. (Ed.) (1987). *The Robot's Dilemma*. Norwood, NJ: Ablex.
- 1831 Riesbeck, C. K. and Schank, R. C. (1989). *Inside Case-Based Reasoning*. Hillsdale, NJ:
1832 Erlbaum.
- 1833 Rochat, P. (2012). Primordial sense of embodied self-unity. In: V. Slaughter and C. A.
1834 Brownell (Eds), *Early Development of Body Representations*. Cambridge, UK: Cambridge
1835 University Press (pp. 3-18).
- 1836 Rosen, R. (1986). On information and complexity. In J. L. Casti and A. Karlqvist (Eds.),
1837 *Complexity, Language, and Life: Mathematical Approaches*. Berlin: Springer (pp. 174-
1838 196).
- 1839 Rubino, G., Rozema, L. A., Feix, A., Araújo, M., Zeuner, J. M., Procopio, L. M., Brukner,
1840 Č. and Walter, P. (2016). Experimental verification of an indefinite causal order. Preprint
1841 arxiv:1608.01683v2 [quant-ph].
- 1842 Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*
1843 61, 85-117.
- 1844 Scholl, B. J. (2007). Object persistence in philosophy and psychology. *Mind and Language*
1845 22, 563-591.
- 1846 Scholl, B. J. and Gao, T. (2013). Perceiving animacy and intentionality: Visual processing
1847 or higher-level judgment?. In M. D. Rutherford and V. A. Kuhlmeier (Eds.) *Social Per-*
1848 *ception: Detection and Interpretation of Animacy, Agency and Intention*. Cambridge, MA:
1849 MIT Press (pp. 197-230).

- 1850 Scholl, B. J., and Tremoulet, P. (2000). Perceptual causality and animacy. *Trends in*
1851 *Cognitive Science* 4, 299309.
- 1852 Seth, A. K. (2013). Interoceptive inference, emotion, and the embodied self. *Trends in*
1853 *Cognitive Sciences* 17, 565-573.
- 1854 Shalm, L. K., Meyer-Scott, E., Christensen, B. G. et al. (2015). A strong loophole-free test
1855 of local realism. *Physical Review Letters* 115, 250402.
- 1856 Shipp, S., Adams, R. A. and Friston, K. J. (2013). Reflections on agranular architecture:
1857 Predictive coding in the motor cortex. *Trends in Neuroscience* 36, 706-716.
- 1858 Simons, D. J. and Ambinder, M. S. (2005). Change blindness: Theory and consequences.
1859 *Current Directions in Psychological Science* 14(1), 44-48.
- 1860 Simons, J. S., Henson, R. N. A., Gilbert, S. J. and Fletcher, P. C. (2008). Separable
1861 forms of reality monitoring supported by anterior prefrontal cortex. *Journal of Cognitive*
1862 *Neuroscience* 20, 447-457.
- 1863 Smith, J. E. and Nair, R. (2005). The architecture of virtual machines. *IEEE Computer*
1864 38(5), 32-38.
- 1865 Steptoe, A., Shankar, A., Demakakos, P. and Wardle, J. (2013). Social isolation, loneliness,
1866 and all-cause mortality in older men and women. *Proceedings of the National Academy of*
1867 *Sciences USA* 110, 5797-5801.
- 1868 Tanenbaum, A. S. (1976). *Structured Computer Organization*. Upper Saddle River, NJ:
1869 Prentice Hall.
- 1870 Tegmark, M. (2015). Consciousness as a state of matter. *Chaos, Solitons & Fractals* 76,
1871 238-270.
- 1872 Toffoli, T. (1980). Reversible computing. In: J. W. de Bakker and J. van Leeuwen (Eds)
1873 *Automata, Languages and Programming: Lecture Notes in Computer Science, Vol. 85*.
1874 Berlin: Springer. pp. 632644.
- 1875 Tononi, G. and Koch, C. (2015). Consciousness: Here, there and everywhere? *Philosophical*
1876 *Transactions of the Royal Society B* 370, 20140167.
- 1877 Trivers, R. L. (2011). *The Folly of Fools*. New York: Basic Books.
- 1878 Turing, A. R. (1936). On computable numbers, with an application to the Entschei-
1879 dungsproblem. *Proceedings of the London Mathematical Society, Series 2* 442, 230-265.
- 1880 van den Heuvel, M. P. and Sporns, O. (2011). Rich-club organization of the human con-
1881 nectome. *Journal of Neuroscience* 31, 15775-15786.
- 1882 Vogel, E. K., Woodman, G. F. and Luck, S. J. (2006). The time course of consolidation
1883 in visual working memory. *Journal of Experimental Psychology: Human Perception and*
1884 *Performance* 32, 1436-1451.
- 1885 von Hofsten, C. (2007). Action in development. *Developmental Science* 10, 54-60.

- 1886 von Uexküll, J. (1957). A stroll through the worlds of animals and men. In: C. Schiller
1887 (Ed.) *Instinctive Behavior*. New York: van Nostrand Reinhold (pp. 5-80). Also published
1888 in *Semiotica* 89 (1992) 319-391.
- 1889 Vossel, S., Geng, J. J. and Fink, G. R. (2014). Dorsal and ventral attention systems :
1890 Distinct Neural Circuits but collaborative roles. *The Neuroscientist* 20, 150-159.
- 1891 Wang, Q., Schoenlein, R. W., Peteanu, L. A., Mathies, R. A. and Shank, C. V. (1994).
1892 Vibrationally coherent photochemistry in the femtosecond primary event of vision. *Science*
1893 266, 422-424.
- 1894 Watson, T. L., Robbins, R. A. and Best, C. T. (2014). Infant perceptual development
1895 for faces and spoken words: An integrated approach. *Developmental Psychobiology* 56,
1896 1454-1481.
- 1897 Watts, D. J and Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks.
1898 *Nature* 393, 440-442.
- 1899 Wiseman, H. (2015). Quantum physics: Death by experiment for local realism. *Nature*
1900 526, 649-650.