

Article

Can A Robot Have Free Will?

Keith Farnsworth

School of Biological Sciences, Queen's University Belfast; Northern Ireland, UK; k.farnsworth@qub.ac.uk

1 **Abstract:** Using insights from cybernetics and an information-based understanding of biological
2 systems, a precise, scientifically inspired, definition of free-will is offered and the essential
3 requirements for an agent to possess it in principle are set out. These are: a) there must be a self
4 to self-determine; b) there must be a non-zero probability of more than one option being enacted;
5 c) there must be an internal means of choosing among options (which is not merely random,
6 since randomness is not a choice). For (a) to be fulfilled, the agent of self-determination must be
7 organisationally closed (a 'Kantian whole'). For (c) to be fulfilled: d) options must be generated from
8 an internal model of the self which can calculate future states contingent on possible responses; e)
9 choosing among these options requires their evaluation using an internally generated goal defined
10 on an objective function representing the overall 'master function' of the agent and f) for 'deep
11 free-will', at least two nested levels of choice and goal (d-e) must be enacted by the agent. The agent
12 must also be able to enact its choice in physical reality. The only systems known to meet all these
13 criteria are living organisms, not just humans, but a wide range of organisms. The main impediment
14 to free-will in present-day artificial robots, is their lack of being a Kantian whole. Consciousness does
15 not seem to be a requirement and the minimum complexity for a free-will system may be quite low
16 and include relatively simple life-forms that are at least able to learn.

17 **Keywords:** self-organization; downward causation; autocatalytic set; goal-oriented behaviour;
18 autopoiesis; biological computing

19 1. Introduction

20 Why do things do what they do? We have a hierarchy of explanations that roughly reflects
21 a gradient in complexity, matched by the epistemic hierarchy which starts with the physics of
22 Hamiltonian mechanics (and Schrödinger's equation), extends through statistical mechanics and
23 complex systems theory, but then declines in power as we try to account for the behaviour of
24 living systems and finally of the human condition, for which we have no satisfactory scientific
25 explanation. One of the most persistent open questions, at the far end of the complexity gradient,
26 is whether we as humans have free will. Here, I attempt to address this question with respect to a
27 broader category of active agent (*sensu* Sharov [1]): anything that can make decisions and act in the
28 physical world. The premise of this paper is that an understanding of the interaction and dynamics
29 among patterns – of the distribution of matter and energy in space and time – may bring such
30 high-level phenomena into resolution. That means a focus on information and its interactions using
31 cybernetics and computation theory, but it also requires a broad concept of information addressing
32 the relationship among patterns (as data) in general, rather than just the statistics of data transmission.
33 If the emergence of material reality (as we experience it) is the assembly of stable configurations (of
34 matter), undergoing transformations and combining as stable composites (e.g. molecules forming
35 materials, forming structures [2]), then a deep understanding of it requires a mathematically precise
36 account of the physics of patterns that are simultaneously the product of material structure and the

37 cause of it [3]. For this investigation, we must focus on the kind of ‘information’ that is embodied by,
38 and processed by natural systems. This concept of intrinsic structural ‘information’ (surely it ought to
39 have a word of its own), is not necessarily ontological (as in the theories of Weizsäcker [4] and Stonier
40 [5]), but refers to at least observable patterns having observable effects [6], is objectively quantifiable
41 [7] and useful in understanding biological processes in terms of cybernetic systems [1], and functions
42 [8,9]. To avoid ambiguity (and conflict) I will refer to this kind of ‘information’ as pattern and the
43 ‘information’ which reduces ‘uncertainty’ in a receiver as Shannon information.

44 I will start by defining what I mean by ‘free-will’ and then give a very brief overview of the
45 philosophical debate about free-will, identifying some of the problems. Then I will introduce ways of
46 thinking about these problems based on cybernetic/computation theory and use these to identify the
47 necessary and sufficient conditions for free-will according to my definition. The class of systems for
48 which these conditions are fulfilled will then be identified and with this, the minimum complexity
49 compatible with autonomous action will be implied. The term ‘robot’ is specifically used here to
50 mean any cybernetic system coupled to a physical system that allows it to independently act in the
51 physical world; this includes living systems (see e.g. the requirements specified for a molecular
52 robot in Hagiya *et al.* [10]) and artificial systems. I am going to conclude that only certain living
53 systems (including humans, together with their anthropic extensions) so far achieve the requirements
54 for organisational autonomy and therefore free-will.

55 1.1. A definition of free-will

56 Free-will is defined here as the condition in which all of the following are jointly true: a) there
57 exists a definite entity to which free-will may (or may not) be attributed, b) there are viable alternative
58 actions for the entity to select from, c) it is not constrained in the exercising of two or more of the
59 alternatives, d) its ‘will’ is generated by non-random process internal to it, e) in similar circumstances,
60 it may act otherwise according to a different internally generated ‘will’. In this definition the term
61 ‘will’ means an intentional plan which is jointly determined by a ‘goal’ and information about the
62 state (including its history) of the entity (internal) and (usually, but not necessarily) its environment.
63 The term ‘goal’ here means a definite objective that is set and maintained internally. The terms used
64 here will be explained and justified in what follows.

65 1.2. The Philosophical background

66 In philosophy, freedom ‘to act as one wills’ is often referred to as ‘superficial freedom’ and the
67 freedom ‘to determine what one wills’ as ‘deep freedom’ (Kane [11] provides a good introduction to
68 the subject, Westen [12] gives a deeper criticism of it). For example, one may be free to drink a bottle of
69 whisky in one sitting and as an alcoholic, or in an irresponsible mood, or emotional turmoil one may
70 will it, but knowing the consequences, one may master the desire and will otherwise: rejecting the
71 opportunity of this poisonous pleasure. If the alcoholism had taken over, or one was under duress,
72 the ‘freedom to choose otherwise’ might be denied and deep freedom would be lost with it, although
73 the superficial freedom to drink would remain. The classical philosophical argument on free-will
74 consists of a) whether it is compatible or not with determinism and b) whether at least some agents
75 (usually people) are at least in part the ultimate cause of their actions.

76 Determinism is the idea that there is, at any instant, exactly one physically possible future
77 [12,13], summarised in the slogan ‘same past: same future’ (see List [14] for a more rigorous
78 analysis). Cybernetics captures determinism in the definition of a determinate machine (DM) as
79 a series of closed, single valued transformations (for example describing a finite state automaton
80 (FSA)). Superficial freedom is often seen as the absence of constraint, leading to the (relatively trivial)
81 conclusion that it is compatible with determinism. But deep freedom needs more than an absence
82 of constraints, it requires alternative paths into the future to provide the ‘freedom to do otherwise’
83 [13]. The cybernetic model of a system with this capacity is of course the non-determinate machine
84 (NDM). However, the NDM is usually conceived as a probabilistic process in which a set of possible

85 states $S\{s_1...s_n\}$ of the system, given the present conditions C (in general including the previous
86 history), may occur at random, with probability set $P\{p_1...p_n\}$ where $\sum p_i = 1$ and each p_i is the
87 probability of each possible state s_i . Most philosophers agree that randomness is not compatible with
88 self-determination, indeed it seems to be the opposite, so they reject random spontaneity as a means of
89 achieving deep freedom (and so do I). They reason that if an agent's action were ultimately caused by
90 e.g. a quantum fluctuation or thermal noise, then we could not reasonably hold the agent responsible
91 for it. This indicates that philosophers supporting the existence of deep free-will are searching for an
92 ultimate cause of actions within the agent of those actions.

93 Unfortunately for them, a paradox arises: since any agent is the product of its composition and
94 previous experiences (the making of it) and these are beyond its control. If it did not make itself and
95 select its own experiences, then its behaviour must be determined by things other than itself. An
96 agent is not free in the deep sense unless it has control over all the events that led to it's choice of
97 action. Recognising that all events in the universe belong to a chain of cause and effect that extends
98 back before the existence of the agent, some philosophers conclude that either a) this deep freedom
99 cannot exist and is considered an illusion (e.g. Van Inwagen [13], reiterated in [15]), or b) the agent is
100 indeterminate so that we get 'same past: different futures'. If they also rule out randomness, then (b)
101 suggests that an agent which could act in more than one possible way from exactly the same state and
102 history (i.e. it is indeterminate) must act without cause. They conclude that this is self-contradictory,
103 hence deep freedom cannot exist. This line of thinking is closely related to Strawson's [16] 'Basic
104 Argument' against free-will, which starts from the premise that for an agent to have free-will it must
105 be the cause of itself and shows, via infinite regress, that this is not possible. Axiomatic to these
106 positions is that the action of an agent can only be either a) random, b) exogenous or c) of itself, with
107 only c) being compatible with free-will.

108 2. Systems with Identity: the closure condition

109 If I am talking of my own free-will, what exactly am I? Only by answering this question do
110 I reach the position of being able to determine if I have free-will or not (rarely, if ever, has this
111 question been posed in the philosophy of free-will). Free-will requires a definite boundary between
112 the internal and external, not only physically (as supplied by e.g. a casing or skin), but more
113 profoundly in organisation and control. For example a computer controlled robot must have all
114 the necessary provisions for physical independence (as in the extraterrestrial exploration robots),
115 but this still leaves them organisationally linked to humanity because their existence is entirely
116 dependent on our gathering and processing the materials for their 'bodies' and assembling these,
117 implicitly embodying them with functional information [1] and programming their control computer
118 (including the goal for operation). For these reasons, they remain extensions of ourselves: tools just
119 as sophisticated hammers. In general, for free-will, the control information of an agent must be
120 independent of anything beyond a cybernetically meaningful boundary. Put the other way round,
121 for the identification of free-will, we must first identify the boundary of the agent, which is defined
122 by independence of control. Given this, the Mars Rover coupled with its human design team seems
123 to meet the closure condition, but the Mars Rover alone does not.

124 This idea of a boundary surrounding a system, such that whatever is within the boundary
125 has the property of organisational independence from what lies without, was encapsulated by the
126 concept of the 'Kantian whole' by Kauffman [17] and can be formally described in cybernetic terms
127 as *organisational closure*.

128 2.1. The Kantian whole

129 A system composed of parts, each of whose existence depends on that of the whole system
130 is here termed a 'Kantian whole', the archetypal example being a bacterial cell [18]. The origin
131 of this terminology lies in Emanuel Kant's definition of an organised whole [19]. For the present
132 purpose, the closure condition is best explained in terms of self-construction, since this implies the

embodiment of self with the pattern-information that will then produce the agents behaviour. In other words, we are to consider a cybernetic system that, by constructing itself materially, determines its transition rules, by and for itself (material self-construction may not be the only way to ensure this self-determination, but assuming the cybernetic relations it embodies are, we may proceed without loss of generality). An autocatalytic chemical reaction network with organisational closure (and this is also what Kauffman [17] considered a Kantian whole) is an anabolic system able to construct itself [20].

Hordijk and Steel [20] and Hordijk *et al.* [21] define their chemical reaction system by a tuple $Q = \{X, \mathcal{R}, C\}$, in which X is a set of molecular types, \mathcal{R} a set of reactions and C a set of catalytic relations specifying which molecular types catalyse each member of \mathcal{R} . The system is also provided with a set of resource molecules $F \subseteq X$, freely available in the environment, to serve as raw materials for anabolism (noting that whilst we are defining an organisational closure, we may (and indeed must) permit the system to be materially and thermodynamically open). The autocatalytic set is that subset of reactions $\mathcal{R}' \subseteq \mathcal{R}$, strictly involving the subset $X' \subseteq X$, which is:

- *reflexively autocatalytic*: every reaction $r \in \mathcal{R}'$ is catalysed by at least one of molecular type $x' \in X'$ and
- composed of F by \mathcal{R}' : all members of X' are created by the actions of \mathcal{R}' on $F \cup X'$.

This definition of an autocatalytic set is an application of relational closure (in computation theory terms, the molecule types are symbol strings for which the Kleene closure is readily defined) and it has been implemented in experiments for exploring aspects of the origin of life (e.g. the GARD system simulating 'lipid world' [22]). Clearly with the two conditions above met, everything in the system is made by the system, but there is a more important consequence. The system is made from the parts (only) and can only exist if they do. Organisational closure of this kind has been identified as a general property of individual organisms [23], many biochemical sub-systems of life [17] and embryonic development [24]).

As it is defined above, living systems fulfil the closure condition, but can we conceive of a non-living system also reaching this milestone? Von Neumann's [25] self-replicating automata show that some purely informational (algorithm) systems have the capacity to reproduce within their non-material domain, but they cannot yet assemble the material parts necessary, nor can they build themselves from basic algorithmic components (they rely on a human programmer to make the first copy). What is needed for the physical implementation of a Kantian whole is the ability to 'boot-strap' from the assembly of simple physical components to reach the point of autonomous replication (i.e. the system must be autopoietic [23,26]). This is necessary to answer Strawson's [16] 'Basic Argument': that for deep free-will an entity must be responsible for shaping its own form and it provides a motivation for rejecting dualism (the idea that the 'mind' is not created from the material universe).

2.2. Emergence and downward causation

Considering the forgoing, we might ask what it is that is making an autopoietic system (e.g. an organism) - is it the components, or the system itself, and in any case, what really is the 'system'. Cybernetics provides an answer to the second question, in that the system is the organisational pattern-information embodied in a particular configuration of interactions among the component parts. Because it is abstract of its material embodiment, it is 'multiply realisable', i.e. composed of members of functionally equivalent parts (see Auletta *et al.* [27]; and Jaeger and Calkins [28] for biological examples). It is not the identity of the components that matters, rather it is the functions they perform (e.g. a digital computer may be 'embodied by' semiconductor junctions, or water pipes and mechanical valves, without changing its identity). Crucially, 'function' is defined by a relationship between a component and the system of which it is a part. According to Cummins [29], 'function' is an objective account of the contribution made by a system's component to the 'capacity' of the whole system. At least one process performed by the component/s is necessary for a process

181 performed by the whole system. This implies that the function of a component is predicated on the
182 function of the whole. This definition was recently modified to more precisely specify the meaning of
183 'capacity' and of whole system, thus: "A function is a process enacted by a system A at organisational
184 level L which influences one or more processes of a system B at level $L + 1$, of which A is a component
185 part" [30].

186 In this context, organisational level means a structure of organisation that is categorically
187 different from those above and below in the hierarchy because it embodies novel functional
188 information (levels may be ontological or merely epistemic in meaning: that is an open debate
189 in philosophy). The self-organisation of modular hierarchy has been described as a form of
190 symmetry-breaking phase transition [31], so the categories either side are quantitatively and
191 qualitatively different. Organisational levels were defined precisely in terms of meshing between
192 macro and micro dynamics (from partitioning the state-space of a dynamic system) by Butterfield [32]
193 and also using category theory to specify supervenience relations and multiple-realizability among
194 levels by List [33]. Neither definition, though, deals specifically with the phenomenon of new
195 pattern-information 'emerging' from the organisation of level L components at level $L + 1$, which
196 is responsible for the emergence of 'new phenomena'.

197 Ellis [34] shows that a multiply realisable network of functions, self-organised into a functional
198 whole, emerges to (apparently) exercise 'downward causation' upon its component parts [34,35].
199 The organisational structure is selecting components from which to construct itself, even though it
200 is materially composed of only the selected components. Since it is purely cybernetic (informational)
201 in nature, the downward control is by pattern-information [28,34] which transcends the components
202 from which it is composed. The pattern-information arises from, and is embodied by, the interactions
203 among the components, and for these reasons it was termed a 'transcendent complex' by Farnsworth
204 *et al.* [36]. Examples are to be found in embryonic development, where a growing cluster of cells
205 self-organises using environmental signals created by the cells taking part [37] and the collective
206 decision making of self-organising swarms (e.g. honey bees in which the hive acts as a unity [38]).

207 There is something significant here for those who conflate determinism with causation. All
208 causal paths traced back would be expected to lead to the early universe. Despite the appearance
209 of near maximum entropy from the uniformity of background microwave radiation, there is broad
210 agreement that the entropy of the early universe was low and its embodied pattern-information
211 (complexity) could not account for the present complexity, including living systems [39]. Novel
212 pattern-information has been introduced by selection processes, especially in living systems, for
213 which Adami *et al.* [40] draw the analogy with Maxwell's demon. Selection is equivalent to pattern
214 matching, i.e. correlation, and is accompanied by an increase of information. Since its beginning, the
215 entropy of the universe has been increasing [39] and some of this has been used as a raw material for
216 transformation into pattern-information. This is achieved by creating the 'order' of spatial correlation
217 through physical self-assembly (atoms into molecules into molecular networks into living systems).
218 This self-assembly embodies new information in the pattern of a higher level structure through
219 the mutual provision of context among the component parts [3]. The process of self-assembly is
220 autonomous and follows a boot-strap dynamic, so it provides a basis for answering Strawson's [16]
221 'Basic Argument' in which the putative agent of free-will is an informational (pattern) structure of
222 self-assembly.

223 2.3. Purpose and will

224 Much of the literature on downward causation uses the idea of 'purpose', though many are
225 uncomfortable with its teleological implication. The aim of this section is to form a non-teleological
226 account of purpose and its connection with will in non-human agents.

227 Cause creates correlation (usually, but not necessarily, in a time-series): the pattern of any action
228 having a cause is correlated with its cause. An action without cause is uncorrelated with anything in
229 the universe and accordingly considered random. If an action is fully constrained, then its cause is the

230 constraint. Thus, freedom from at least one constraint allows the cause to be one of either: random,
 231 or exogenous control or agent control (in which ‘control’ means non-random cause). By definition the
 232 cause is only taken to be the agent’s will if it originated in agent control. Correlation alone, between
 233 some outcome variable x and some attribute a of the agent, is not sufficient to establish will: a) because
 234 correlation has no direction (but metrics such as ‘integrated information’ [41] can resolve direction)
 235 and b) because a may itself be random in origin and thereby not of the agent’s making. Marshall *et al.*
 236 [42] showed that cause can be established at the ‘macro’ level of agent (as opposed to the ‘micro’ level
 237 of its components) using an elaboration of integrated information, so the pattern in x can be attributed
 238 to agent-cause. Because the agent-based cause could be random ((b) above), we must form and test
 239 a hypothesis about the effect of x on the agent before we can attribute the cause to the alternative of
 240 agent-will. The hypothesis is that the effect of a on x is to increase the overall functioning F of the
 241 agent. If this were true, then to act wilfully is to reduce the entropy of x , by increasing the probability
 242 of an outcome x' where $x' \Rightarrow F' > \bar{F}$ (and \bar{F} is the average F). That means that the mutual information
 243 between a wilful action $a(t)$ at time t and the resulting function $F(t + \tau)$, $\tau \geq 0$ is greater than zero.
 244 This mutual information between action and future functioning is taken to imply a ‘purpose’ for the
 245 action, so purpose is identified by the observations that:

$$H(a) + H(F) > H(a|F) \text{ and } F|a > F|r, \quad (1)$$

246 where r is a random (comparator) variable. This is clearly an observational definition and is
 247 in some way analogous to the Turing test, but it is a test for purpose rather than ‘intelligence’. It
 248 represents our intuition that if a behaviour repeatedly produces an objectively beneficial outcome for
 249 its actor, then it is probably deliberate (repeatedly harmful behaviour is also possibly deliberate, but
 250 all such actions are regarded as pathological and thereby a subject beyond the present scope).

251 To recap, for attributing the action to the agent’s will, we must at least identify a purpose so that
 252 Eq. 1 is true. The purpose is a pattern embodied in the agent, which acts as a template for actions of
 253 the agent that cause a change in future states (of the agent, its environment, or both). We may call this
 254 pattern a ‘plan’ to attain an objective that has been previously set, where the objective is some future
 255 state to which the plan directs action. Specifically, let the objective be a state X (of the system or the
 256 world, etc.), which can be arrived at through a process P from the current state Y , then the purpose is
 257 a ‘plan’ to transform $Y \rightarrow X$ by the effect of at least one P and at least one function F is necessary for
 258 the process P to complete.

259 The homeostatic response to a perturbation, for example, has maintenance at the set-point as
 260 its purpose. Y is the perturbed state, F is some function of the internal system having the effect of
 261 causing a process P , i.e. some transition $Y \rightarrow X$. In general there is more than one P and more
 262 than one F for achieving each. This results in a choice of which to use: it is a choice for the agent
 263 described by the system. To make a choice requires a criterion for choosing (else the outcome is
 264 random and therefore not a choice). The criterion for choosing is a ‘goal’ G , consisting of one or
 265 more rules, which identify a location in a function describing the outcome (which we may call the
 266 *objective function*). In general, this location could be any and it is essential for freedom of will that it be
 267 determined by the agent of action alone. However, in practice it is most likely to be an optimisation
 268 point (in living systems, this is implied by Darwinian evolution and in designed systems, it is the
 269 basis of rational design). So, narrowing the scope, but with justification, let us take the criterion
 270 for choosing to be a ‘goal’ G , consisting of one or more optimisation rules. For example, of all the
 271 possible systems performing homeostasis, the purposeful one is defined as enacting P' such that
 272 $Y \rightarrow X$, with $P' \in P | \max(\mathcal{G})$, where \mathcal{G} is an *objective function* for which the optimisation goal $G(\mathcal{G})$
 273 is satisfied, contingent upon the options (e.g. P proceeds as quickly as possible, or with minimum
 274 energy expenditure, etc.). Accordingly, ‘will’ is defined by a purpose which is a plan to enact a process
 275 causing a transition in state, ‘as well as possible’ (according to $G(\mathcal{G})$) – notwithstanding the earlier
 276 comment about pathological purposes. A free-will agent has a choice of transition and a goal which

277 identifies the most desirable transition and the best way to enact it, from those available. These two
278 choices can be united (by intersection), without loss of generality, to one choice of best transition.

279 One of the reasons for objecting to teleological terms such as 'purpose', 'plan' and 'goal' in
280 relation to natural systems has been the belief that a plan implies a 'designer' - the concept at the
281 centre of the most famous battles between science and religion. This implication is not necessary and
282 is rejected here (following the argument of Mayr [43]). A plan is merely a pre-set program of steps
283 taking the system from Y to X; it is the concept for which computation theory was developed. It may
284 be designed (the work of an engineer), but also may have evolved by natural selection (which also
285 supplies the goal - in which case it is a teleonomic system (*sensu* Mayr [43])).

286 A plan, as an ordered sequence of transformations, is an abstraction of information from the
287 physical system, which for free-will must be embodied within the system. A more subtle implication
288 of 'plan' is that as a path leading from Y to X, it is one among several possible paths: different plans
289 may be possible, perhaps leading to different outcomes. There is a fundamental difference between
290 this and the inevitability of a dynamic system which follows the only path it may, other than by the
291 introduction of randomness. The reason is that for a dynamical system all the information defining
292 its trajectory is pre-determined in the initial (including boundary) conditions and the laws of physics.
293 The initial conditions constitute its one and only 'plan'. If a system embodies pattern-information
294 (by its structure) which constitutes a developed plan, then this pattern-information may direct the
295 dynamics of the system along a path other than that set by the exogenous initial conditions (though
296 we may consider the structure of the system to be a kind of initial condition). The point is that the
297 embodied plan gives freedom to the system, since it 'might be otherwise' - there could be a different
298 plan and a different outcome. We see this in the variety of life-forms: each follows its own algorithm
299 of development, life-history and behaviour at the level of the individual organism. The existence of a
300 plan as abstract pattern-information is a pre-requisite for options and therefore freedom of action.

301 2.4. Goal, master function and will nestedness

302 Now let us complete the connections between will, goal and function. The previous argument
303 reveals an important difference between downward and any other kind of causation (considered
304 important by Walker [44]): the former must always be directed by a purpose, for which we need to
305 identify a goal (upward and same level causation are satisfactorily explained by initial conditions
306 [34]). Viewing entities and actions both as the consequence of information constraining (filtering)
307 entropic systems, then the role that is taken by initial conditions in upward causation, is taken by
308 system-level pattern-information (the transcendent complex [36]) in downward causation. Since the
309 goal G is a fixed point in an objective function \mathcal{G} , it constitutes information (e.g. a homeostatic
310 set-point) that must be embodied in the agent's internal organisation. Since the objective function
311 \mathcal{G} represents the overall functioning of the system (at its highest level), it matches the definition given
312 by Cummins [29] and Farnsworth *et al.* [30]. The highest level function from which we identify the
313 purpose of a system was termed the 'master function' by Jaeger and Calkins [28], so the will of an
314 agent is instantiated in the master function. This then identifies \mathcal{G} with 'master function' and 'will'
315 with $G(\mathcal{G})$.

316 Ellis [34] identifies 5 types of downward causation, the second being 'non-adaptive information
317 control', where he says "higher level entities influence lower level entities so as to attain specific fixed
318 goals through the existence of feedback control loops..." in which "the outcome is not determined
319 by the boundary or initial conditions; rather it is determined by the goals". Butterfield [32] gives
320 a more mathematically precise account of this, but without elaborating on the meaning or origin
321 of 'goals'. Indeed, as both Ellis [34] and Butterfield [32] proceed with the third type: downward
322 causation "via adaptive selection" they refer to fitness criteria as "meta-goals" and it is clear that these
323 originate before and beyond the existence of the agent in question. Ellis [34] describes meta-goal as
324 "the higher level 'purpose' that guides the dynamics" and explains that "the goals are established
325 through the process of natural selection and genetically embodied, in the case of biological systems,

326 or are embodied via the engineering design and subsequent user choice, in the case of manufactured
327 systems".

328 This suggests a nested hierarchy of goal-driven systems and for each, the 'goal' is the source
329 of *causal power* and as such may be identified as the 'will' (free or otherwise). We may interpret the
330 definition of 'deep freedom' [11] as meaning that an agent has at least two nested levels of causal
331 power, the higher of which, at least, is embodied within the agent (as causal pattern-information).
332 This concept may be formalised after introducing the discrete variable 'will-nestedness' \mathcal{N} which
333 counts the number of levels of causal power exercised over a system, *from within the agent as a whole*
334 (i.e. at the level of master function), the \mathcal{N} th level being the highest-level internal cause of its actions.

335 Among organisms in general, the master function specifies the criteria by which the organism
336 is to assess its possible future reactions to the environment. It is so much an integral part of the
337 organism that without it, the organism would not exist. However, it was not chosen by the organism
338 (in the sense of deep free-will) because it was created by evolutionary filtering and inherited from
339 its parent(s) - as all known life has been created by the previous generation copying itself. For
340 single celled organisms the biological master function is to maximise their cell count by survival
341 and reproduction, but in multicellular organisms, this master function exists, by definition, at
342 the level of the whole organism (the unconstrained drive to proliferate a single cell line leads to
343 cancer). Organisms with a central nervous system, regulated by neuro-hormone systems, with their
344 corresponding emotions, can implement more complex (information rich) and adaptable (internally
345 branched) algorithms for the master function, which may include will-nestedness $\mathcal{N} > 1$. In humans,
346 this is taken to such an extent that the biological master function may seem to have been superseded
347 (but the weight of socio-biological evidence may suggest otherwise [45,46]).

348 2.5. The possibility of choice and alternative futures

349 So far I have identified organisational-closure and the internal generation of a goal-based plan
350 as prerequisites for free-will, but have not yet addressed the strict possibility of freedom for an
351 agent constructed from elemental components that necessarily obey physical determinism. List [14]
352 provides a philosophical argument for meeting this requirement, constructed from supervenience and
353 multiple realisation of an agent in relation to its underlying (micro) physical level: "an agent-state is
354 consistent with every sequence of events that is supported by at least one of its physical realizations"
355 [14]. He shows that this may apply not only to multiple micro-histories up to t , but in principle
356 includes subsequent $(t + \tau)$ sequences at the micro-level, which may map to different agent states
357 and therefore permit different courses of action at the (macro) agent-level. To explain: for any given
358 time t , the macro-state $Q_i(t)$ is consistent with a set of micro-states \mathbf{s} , at least one of which $s_i \in \mathbf{s}$
359 may lead (deterministically) to a new state $s_j \in \mathbf{s}$ at $t + 1$, with which a different macro-state Q_j is
360 consistent, thus giving the agent a choice of which micro-state history to 'ride' into $t + 1$ (this idea is
361 developed with rigour by List [14], and illustrated with 'real-world' examples; it is the basis on which
362 he concludes that agents may be 'free to do otherwise', despite supervening on deterministic physical
363 processes).

364 The possibility of choice may exist at multiple levels of system organisation within a hierarchical
365 structure, re-applying the same principles as identified by List [14] for each level of macro-micro
366 relations. For any system level L to have this potential, it must have the attributes of an agent-level:
367 supervenience and multiple realisation such that pattern-information with causal power emerges at
368 level L from $L - 1$: i.e. a transcendent complex exists at level L . But this does not necessarily give *free*
369 *will* to a system of that level, since for that, it must be organisationally closed. If it were not so, we
370 would not be able to identify the system at level L as an entity to which free-will could be ascribed.
371 Thus will-nestedness cannot be attributed to levels of organisation below that of the Kantian whole.
372 Since the Kantian whole is, by definition, the highest level of organisation to which free-will may
373 be ascribed (any causal power beyond it rules out its free will), then will-nestedness can only apply
374 at the level of the Kantian whole. Given this, the will-nestedness must be constructed from purely

375 organisational, i.e. pattern-informational and therefore be purely computational in nature. This is an
376 important deduction: free-will can only be an attribute of a Kantian whole and it can only result from
377 the cybernetic structure at the level of the Kantian whole.

378 3. Choosing possible futures: The *computational condition*

379 We see that for free-will, an agent must have an independent and internally generated purpose
380 for action and that this requires it to be organisationally closed. Free-will further requires the agent
381 to use this purpose to choose among options. To do so, it needs an internal representation of *possible*
382 futures from which to choose and an internally generated means of choosing. We now turn to the
383 conditions which enable these essentially computational facilities.

384 3.1. Information abstraction

385 Information abstraction at the organisational boundary is crucial to achieving autonomy. It strips
386 off the physical effect of the external environment to take only the abstract information as a signal.
387 Causes are transformed into signals, their effects being rendered responses (which thereby may
388 become optional). It is this separation of information (as signals) from the physical force of cause and
389 effect that releases the agent from attachment to the cause-effect determinism of its environment. The
390 material apparatus for performing this task is a transducer: the tegumental membranes of bacteria
391 and other cells contain a wide variety of transducers (receptors) and we expect an artificial robot to
392 be well equipped with them too. This is not a merely technical point. The closure condition gives
393 the system a degree of causal independence and the boundary transducers give it a sensitivity to its
394 environment whilst preserving this independence. The boundary is the place where the inevitability
395 of cause and effect of the environment meets that of the internal processes of the system and the
396 transducers are the interface between these causal chains. Internally to the system, the environment
397 is 'reduced' to abstract, representational, information by the transducers [47]. Now the question is,
398 what must the agent do with this information in order to exercise free-will?

399 Of course the answer is to compute: more specifically, to perform transformations on the data
400 as a result of a sequence of physical changes in the physical structure of the agent. Such changes
401 are described by automata theory, for which the most basic automaton has two states and can
402 potentially change state on receiving a signal to which it responds: it is a switch (e.g. a protein
403 molecule with two conformations is an 'acceptor' of all strings from the alphabet $\{0,1\}$, where
404 these symbols may represent the presence / absence of e.g. another molecule or a level above a
405 threshold (e.g. temperature)). Obviously, the switch is the elemental component for generating
406 *discrete options*. A less obvious, but crucial property of the switch for the physical embodiment
407 of computation is 'thermodynamic indifference'. Walker and Davies [48] focus on computation in
408 explaining the origin of life, referring to genetics-first theories as 'digital-first', emphasising the need
409 for 'programmability' and its provision by informational polymers (the genetic oligomers RNA, DNA
410 etc.). By programmability they meant that components of a system are configured so that the system
411 state can change reversibly, approximately independent of energy flow: i.e. changes of state are not
412 accompanied by substantial changes in potential (stored) energy. If they were, then switching would
413 always be biased by the difference in energetic cost between e.g. switching on and off. Energetically
414 unbiased switching is the physical underlying mechanism of information abstraction referred to by
415 Walker and Davies [48]. In reality, switching (and state-changes in general) always have energetic
416 consequences (more deeply, there is always an exchange of entropy between the system and an
417 external energy source), which is one of the reasons an autonomous agent must complete work cycles
418 as Kauffman [49] specifies. What makes the informational polymers (e.g. DNA) of life special is the
419 fact that they are reversible in a way that is thermodynamically indifferent (or very nearly so - see
420 Ptashne [50]). Any ordered set of n switch positions (e.g. 1,1,0,0) has very nearly the same potential
421 energy as any other ordered set of n (e.g. 0,1,0,1).

422 For free-will, both the self-assembly of autocatalytic systems and the computational
 423 requirements (switches and memory) are jointly necessary. Walker and Davies [48] and Walker [44]
 424 proposed that the autonomy of living systems arises from the combination of ‘analogue’ chemical
 425 networks and ‘digital information processing’.

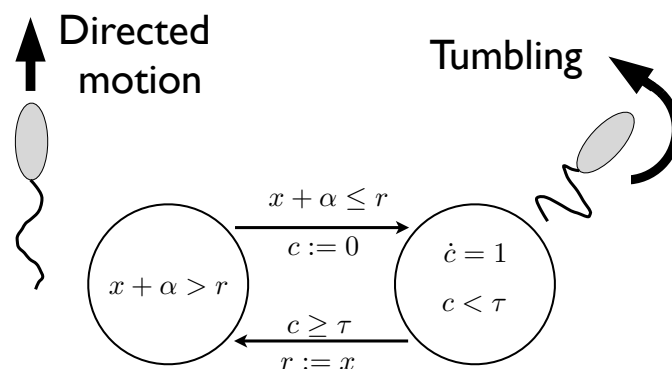


Figure 1. A bacterial chemotaxis controller described as a hybrid automaton, realised in practice by e.g. *E-coli* species, but also as an engineering design for a molecular robot, using DNA-based components by Hagiya *et al.* [10], from which the figure, slightly modified, is taken. Note that α and τ are internal set-points, which constitute pattern-information embodied in the molecular robot’s structure and are ultimately determined by Darwinian evolution (natural robots) or intentional design (human artefacts). The bacterium responds to the environmental concentration (e.g. of glucose) - its objective function x , which is detected by a membrane transducer which generates the internal variable r . It is searching for a higher concentration when x is below a set threshold (directed motion) and tumbles randomly for a time set internally by τ using an internal ‘clock’ signal c , whenever the concentration is at least equal to the threshold. The swimming and tumbling are modes of action of its flagella. The objective function of this system is the experienced concentration x and the goal is $\max(x)$. The chemotaxis controller exhibits will-nestedness of 1.

426 A hybrid automaton is a good model for the construction of a free-will agent. Fig. 1 (from
 427 Hagiya *et al.* [10]) shows an example of a hybrid automaton which combines discrete-state with
 428 continuous dynamical systems, such that the discrete states (as modes of functioning) determine the
 429 system’s responses to dynamic variables and these responses potentially influence the trajectory of
 430 the dynamics. In this pedagogical example from Hagiya *et al.* [10], there are two state variables α and
 431 τ which determine the set points for autonomous chemotaxis behaviour in the bacterium represented.
 432 This system can freely maximise its (experienced) environmental concentration of x (e.g. glucose), so
 433 its goal $G(\mathcal{G})$ is defined, but it cannot choose how (hence $\mathcal{N} = 1$), so it cannot express free-will in the
 434 deep sense. However, if it had independent control over α and τ , with the ability to adjust these values
 435 according to a plan of its own making, then it would fulfil the condition of having willed its own
 436 behaviour ($\mathcal{N} = 2$), at least in the sense defined in the previous section. As part of a living bacterium,
 437 the system depicted would be a component of a Kantian whole (the free-living organism), so all that
 438 is missing for free-will is a plan for determining α and τ according to an internally generated goal (a
 439 master function) and some means of computing this. The computational requirements for free-will
 440 are identified as follows.

441 3.2. The representation of self

442 Firstly, a free-will agent must maintain an internal representation of itself, and also the effect of
 443 its environment on its internal state, to enable it to assess each of its options for action.

444 At the root of automata theory lies an attempt to fully describe (and therefore predict) the
 445 behaviour of a system without a detailed mechanistic account of its internal workings (the black box
 446 approach). The system is captured in the mapping between environmental stimuli and responses:

$$R(t+1) = W(\mathbf{S}(t), S(t)) \quad (2)$$

447 where $R(t)$ and $S(t)$ are the response and stimulus, and $\mathbf{S}(t)$ is the history of stimuli experienced
 448 by the system, *from the beginning of its formation* up to t . This presents an immediate problem, since
 449 in general $\mathbf{S}(t)$ is arbitrary and infinite in range (it is instructive to think of Strawson's [16] 'Basic
 450 Argument' in terms of Eq. 2: the response of a system, in general, depends on its environment from
 451 before the system came into existence). The solution to this indefinite $\mathbf{S}(t)$ problem (provided by
 452 Moore [51]) is to assume that the infinite set of $\mathbf{S}(t)$ may be partitioned into a finite number of disjoint
 453 equivalence sets, each containing the histories that are equivalent in their effect on $R(t+\tau)$, $\forall \tau$, where
 454 τ is in the interval $[0, \infty]$. These equivalence sets are represented as states $Q(t)$ of the system, so that:

$$R(t+1) = W(Q(t), S(t)) \text{ and } Q(t+1) = U(Q(t), S(t)), \quad (3)$$

455 where $U(Q(t), S(t))$ is the transformation function dictating the transition of system state given its
 456 present state and that of the stimulus (see Minsky [52]. p16-17). Thus $Q(t)$ represents how the system
 457 is now, *given* its previous history of experiences. $Q(t)$ corresponds to the agent-state of List [14],
 458 which is multiply-realizable and which is, at least in principle, free to take more than one value at
 459 some point in the future $t+\tau$, despite supervening on a wholly determined set of micro-histories.
 460 List [14] showed the possibility of choice at the agent-level, but this does not necessarily mean that
 461 $Q(t)$ is indeterminate. Specifically, for free-choice (of $Q(t+1)$ and implied $R(t+1)$), the direction
 462 taken at the branching point t must be determined by a process internal to the agent that represents
 463 its 'purpose' (as defined earlier). For the choice to be purposful, it must be based on an assessment of
 464 the outcomes that would arise from choosing each of the options. This entails a prediction of possible
 465 futures, for which a free-will system must have a model of itself in its environmental context.

466 The question now is, how can a system create such a representation by and for itself - not
 467 'programmed' by some exogenous source of information? The answer seems to be as it is with
 468 material self-assembly: a boot-strap, step by step, gathering of pattern-information embodied in form,
 469 such that as the form grows, it increases in complexity. In the particular case of building a model
 470 of self and environment, this process is one of learning, for which the field of 'machine learning'
 471 provides our understanding. Well known advances in this field have already led to sophisticated
 472 learning among pre-existing (i.e. not self-assembled) computation systems such as deep neural
 473 networks etc.. The difference here is that the learning is not merely a statistical problem, but one
 474 of simultaneous self-construction, which must begin with simple systems, so in the remainder of this
 475 section, only basic and simple systems capable of unsupervised learning are discussed.

476 3.2.1. Learning in a constant environment

477 The most basic form of learning is operant conditioning (reinforcement learning), described
 478 mathematically by Zhang [53] as follows. Let R_i be one of a set of N possible responses ($R_i \in \mathbf{R}$, $i \in$
 479 $[1, N]$) in a constant environment, occurring with probability r_i . For each response there is a 'reward'
 480 ρ (which coincides with the objective function \mathcal{G} that defines the goal of the system: $\rho \rightarrow f(\mathcal{G})$), so
 481 that the incremental change in probability of the i th response is:

$$\Delta r_i = a \rho_k (\kappa_{k,i} - r_i), \quad (4)$$

482 where $\kappa_{k,i}$ is the Kronecker delta function (equal to 1 if $i = k$, else equal to zero) and a is a
 483 learning rate constant. Since the average change in response over the ensemble of possible responses
 484 is the frequency-weighted sum: $\Delta \hat{r} = \sum_k^N r_k \Delta r_k$, the result is that the frequency of the i th response

485 incrementally increases in proportion to the difference between its reward and the average over all
486 rewards:

$$\Delta r_i = a r_i (\rho_i - \hat{\rho}), \text{ where } \hat{\rho} = \sum_k^N r_k \rho_k. \quad (5)$$

487 The dynamic quantified in Eq. 5 describes learning by maximising the reward experienced. Such
488 learning is equivalent to making an increasingly accurate model of the (static) relationship between
489 the agent's internal state and the environment, via (Bayesian) 'trial and error' sampling of responses.
490 Given a constant environment, the solution to Eq. 5 yields a single, reward maximising response:
491 $R^* = R_i$ such that $\rho_i = \hat{\rho}$ and $r_i = \kappa_{k,i}$.

492 To achieve this in practice the system must at least keep a record of the reward for the last
493 response made and the average reward, for which an automaton with an external memory is required
494 (e.g. a push-down automaton, though this is still essentially a DFA). Quantifying the complexity of
495 such a system is probably best achieved through a programme (algorithmic) complexity measure
496 since the information instantiated by such an automaton is almost all in its transition mapping
497 and there are robust methods for reducing this to the minimum description, leading directly to
498 the Kolmogorov complexity. The process of learning can be interpreted in information terms: if
499 the starting probability distribution of \mathbf{R} is \tilde{r} , then the initial Shannon entropy of the system is
500 $H = -\sum_i^N \tilde{r}_i \log(\tilde{r}_i)$ and the final entropy is zero: having completed its learning, the system has
501 no uncertainty about the best way to respond to this environment. In this state, the automaton is a
502 complete representation of its interaction with its environment (i.e. the distribution of rewards over
503 its repertoire of responses) and it embodies exactly H units of information: a quantity which should
504 match the algorithmic complexity measure (though not tested here).

505 3.2.2. Extension to a variable environment

506 Generalising to a variable environment, for which a set of finite states \mathbf{S} is an adequate
507 representation, there would exist a reward maximising response for each state: $R_i^* \rightarrow (S_i)$, such
508 that R_i^* solves Eq. 5. The agent may choose to maximise its reward over all \mathbf{S} , and to enable this, it
509 must learn the best response for every $S_i \in \mathbf{S}$. In information terms, first let $H(\mathbf{s})$ be the entropy
510 of the environment having probability distribution \mathbf{s} and $H_t(\mathbf{r}_t)$ be that of the responses, given their
511 probability distribution \mathbf{r}_t at time t . The Shannon information the agent has about its environment (in
512 terms of its rewards) is:

$$\begin{aligned} I_t(\mathbf{s} : \mathbf{r}_t) &= H(\mathbf{s}) + H(\mathbf{r}_t) - H(\mathbf{s}, \mathbf{r}_t) \\ &= H(\mathbf{s}) + H(\mathbf{r}_t) - H(\mathbf{s}|\mathbf{r}_t), \end{aligned} \quad (6)$$

513 meaning that the agent is learning both the distribution \mathbf{s} and the reward associated with each
514 $S_i \in \mathbf{S}$. This *mutual information* is embodied in the structure of the agent and can be used as a measure
515 of its complexity. The structure of the agent may be too simple to embody as much as the maximum
516 mutual information, in which case its learning will be limited and it will not make optimal responses,
517 so Eq. 6 is a measure of the minimum complexity required for optimal behaviour from the agent. It
518 would be possible for an agent to implement this learning system by 'growing' multiple copies of the
519 DFA with memory (one for each $S_i \in \mathbf{S}$) that is used for a constant environment. The output of each
520 of these would then be the input to a further DFA which it uses to maximise the reward across all
521 of them. This 'growth' would be enacting meta-learning: the agent would increase its complexity in
522 response to rewards. The number of states in \mathbf{S} is not known by the agent a-priori, so the number
523 of DFAs needing to be 'grown' is indeterminate. Further, account should be taken of the extended
524 time needed to perform such laborious learning and the consequences of the agent being wrong in
525 so many trials. It seems that for practical reasons there comes a point when a more powerful kind of

526 computation becomes necessary. In computation terms, such a learning problem requires at least a
 527 finite and non-volatile memory, effectively to store multiple instances of the single learning problem
 528 encountered in a constant environment. For this reason a Turing Machine would be a more realistic
 529 option.

530 3.3. The Free-Will Machine

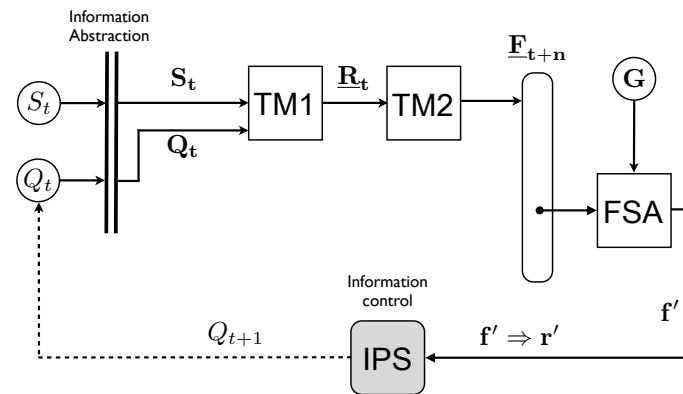


Figure 2. A conceptual ‘free-will’ machine, which generates predictions of its state in alternative futures E_{t+n} using an internal representation of itself interacting with its environment, selecting the optimal from among these, using a goal-based criterion G , in which the goal is internally determined (further explanation in the text).

531 These computational requirements for free-will to become possible are brought together in
 532 the hypothetical ‘free-will machine’ of Fig. 2. This information processing must be implemented
 533 by the agent to which free-will is ascribed and that agent must, further, be a Kantian whole
 534 for the requirements of free-will to met. The current state of the environment (external) and of
 535 the agent (internal) are derived by information abstraction from the physical world: the array of
 536 receptor molecules in the cell membrane, the nervous senses of an animal, or the transducers of a
 537 human artefact all perform this task. The first Turing machine implementation TM1 constitutes a
 538 representation of the agent in the present (relevant aspects of its environment are included) and is
 539 informed (updated) by the state information Q_t and S_t . The function of this representation is to
 540 identify the set of possible responses R_t that the agent can make, given Q_t, S_t (underline notation
 541 denotes a set).

542 TM2 uses these hypothetical responses to compute the set of possible futures at a time $t + n$ (n
 543 may take any positive value) E_t for which in a simple case $TM2 : R_t \rightarrow E_t$ is a map of responses
 544 onto possible futures (simple, because there is not necessarily a 1:1 relation between R_t and E_t , but
 545 we need not be concerned with that at present). There is no limit on the number of possible futures
 546 that TM2 may compute, but it must be at least 2 for a choice to exist. These possible futures are each
 547 represented by a set of states f_n , each member f_n being equivalent to a prediction of a possible Q_{t+n} .
 548 The FSA chooses from among these, using a selection criterion based on the objective function defined
 549 by a goal G , which is generated by the agent (not exogenously). This goal is the maximisation of the
 550 master function, (e.g. for a living agent this is life-time reproductive success). The goal enables the
 551 optimal possible future state to be recognised (it is the one which maximises the master function) and
 552 this future state f' implies an optimal response r' (in general there could be more than one, in which
 553 case the agent will be indifferent among them). Having selected an optimal response, the agent then
 554 must implement it in the physical realm. Since the computation of r' has been conducted in the realm
 555 of information, this step appears to involve the control of material by information. In practice all
 556 the computation and indeed all the information is instantiated by material and energy acting in the

557 physical realm, so our cybernetic model is merely an abstract representation of the organisation of the
558 physical processes which lead to the implementation in the physical system and this return to physical
559 reality is represent by the action (IPS). Such implementation inevitably results in a transformation of
560 the agent into a new state Q_{t+1} , together with S_{t+1} and this restarts the cycle. It may be noted that Von
561 Neumann's self-replicating automata are proven universal Turing machines [25] and Turing machines
562 are thought to be common among living systems, so this computational arrangement is not beyond
563 the bounds of possibility.

564 4. Discussion and Synthesis

565 To summarise, the essential requirements for free will are:

- 566 ● R1) There must be a self to self-determine.
- 567 ● R2) There must be a non-zero probability of more than one option being enacted.
- 568 ● R3) There must be an internal means of choosing among options (which is not merely random,
569 since randomness is not a choice).

570 for R3 to be fulfilled:

- 571 ● R4) Options must be generated from an internal model of the self which can calculate future states
572 contingent on possible responses.
- 573 ● R5) Choosing among these options requires their evaluation using an internally generated goal.
- 574 ● R6) For 'deep free-will', at least two nested levels of choice and goal (R4-R5) must be enacted by
575 the agent.

576 R1 and references to 'internally generated' are fulfilled by organisational closure. For R2, the
577 possibility of options, which implies 'multiple futures' has been established for the level of the agent
578 by List [14]. R3 and its predicates R4-R5 imply a minimum level of computational power, which
579 in principle can be met by a small set of Turing machines, which may in principle be implemented
580 by a Von Neuman architecture computer, a network or cellular automaton-based or any other sort
581 of computer, including a biomolecular system such as found in higher animal life, but it seems to
582 be beyond the power of a single living cell (though that last point is not yet established). R4 in
583 particular seems to require a finite memory (the size depends on the complexity of the agent and
584 its environment) and R6, the qualifier for deep free-will, adds a little more to the computing power
585 necessary, but it is important to note that this extra is not a step-change: it is not qualitatively more
586 demanding than the automated decision making required by R3.

587 The question of free-will is not one of whether an agent's actions are caused, since all actions
588 ultimately have a cause. The ultimate cause of any action can be understood as resulting from
589 selection over random actions by a pattern, which leaves a correlation with the pattern that caused
590 the selection (instantiating pattern-information). All living organisms, including people, were
591 produced by information-pattern filtering, proximally by molecular replication (creating inheritance)
592 and ultimately by Darwinian selection. All human artefacts were created by following a design
593 pattern (though it may not have been completed before artefact construction), so they correlate with
594 their design. Even inanimate objects, such as stars, lakes and sand grains, owe their form to the
595 information-pattern of underlying physical laws, Pauli's exclusion principle and the distribution of
596 matter and energy in space following the big-bang. To this, we must add randomness which has
597 been entering as 'informational raw material' into the universe, disrupting the original patterns and
598 opening opportunities for novelty (evolutionary for life) and more widely directing the course of the
599 universe in unexpected ways as its history tracks a course in the highly ergodic space of possibilities.

600 Taking Strawson's [16] Basic Argument seriously, this pattern-correlation and the injection of
601 randomness both deny free will. From them, we obtain a model in which the identity of all
602 things, including human beings, is an illusion: as if the universe was all one complex manifestation
603 which only *appears* to include separate agents. Closer inspection shows how the nested-hierarchical

604 construction of this complexity entails the creation of genuinely new pattern-information, caused
605 by and embodied in the interactions among component parts of putative agents. This novel
606 pattern-information transcends its component parts and can exert downward causation upon them.
607 Some structures (such as autocatalytic sets) created this way are organisationally closed (though
608 materially and thermodynamically remain open systems). Because of this, their internal dynamics
609 are, at least partially, separated from the external dynamics of their environment and this gives
610 them an organisational boundary, enabling internal to be defined against external. At this boundary,
611 external and internal chains of cause and effect interact through transducers which transform
612 physical determinism into stimulus-response relations. Systems with these properties are essentially
613 cybernetic and although their low-level processes are continuous with the rest of the universe, List
614 [14] has shown that in principle they may have options for their next state and response: they
615 are freed from physical determinism. To translate this freedom into free-will, requires that the
616 (partially) independent agent chooses from among its options and this entails an internal computation
617 of possible futures and their evaluation against a goal representing the fulfilment of the agents'
618 'master function' (i.e. its purpose). This goal is a fixed point in an objective function which may be
619 simple (as in a homeostatic system), but also arbitrarily complex and multi-layered, taking account of
620 multiple time-scales and interactions with other agents. If the objective function is at least two-layered
621 (will-nestedness $\mathcal{N} \geq 2$), then it effectively has a choice of what to choose and thereby could fulfil the
622 established definition of (deep) free-will [11,12]. This calls into question the idea that free-will is an
623 all or nothing capacity, instead, it suggests free-will to be a discrete quantity and even something we
624 could in principle measure as a trait of a system.

625 The reason is that deep free-will has so far been defined as the freedom to choose ones will,
626 but the analysis presented here shows that to be wilful, a choice must be purposeful, which means
627 optimising an objective function. Freedom is in the choice of objective function. Since therefore, the
628 core of will is the objective function, deep freedom is the freedom to choose this, but to be wilful,
629 this choice in turn must optimise a hierarchically superior objective function, which must have been
630 determined by something. We can conceive of a large but finite nested set of such objective functions,
631 but *ultimately* the highest of them all must be provided either arbitrarily (e.g. at random) or by
632 natural selection (or its unnatural equivalent), or by design: in all cases, not the free choice of the
633 agent. This applies even to human beings, who are still subject to 'design' by inheritance, selected
634 by evolution. The depth of free-will is therefore a discrete, finite variable ('will-nestedness' \mathcal{N}) and
635 we may speak of one kind of agent being more deeply free than another, but no kind of agent can
636 be *ultimately* free-willed in the sense meant by deep free-will (presumably, humans attain the highest
637 will-nestedness of known agents).

638 If interpreting the philosophers' definition of deep free-will as $\mathcal{N} \geq 2$ is correct, then it is not
639 hard to achieve in principle. The impediment to a non-life robot acquiring that sort of free-will is
640 not computational, it is the closure constraint with its requirement for bootstrapping self-assembly
641 (especially since this includes the 'growth' of the sort of computational apparatus indicated in Fig. 2).
642 That is a problem already solved by life. Quite likely $\mathcal{N} \geq 2$ in most or all organisms having
643 at least a limbic system, so free-will defined this way can be attributed to most or all vertebrates
644 [54]. The computational requirements for exercising purposeful choices are not very challenging for
645 artificial computers. Among human artefacts, including what most people would define as robots,
646 the organisational closure condition is the major hurdle not yet leapt. For the time-being, it seems
647 free-will, as defined here, is a unique property of living things.

648 **Acknowledgments:** This work was unfunded, but was inspired by attendance at the workshop "Information,
649 Causality and the Origin of Life" in Arizona State University at Tempe, AZ Sept. 30th - Oct. 2nd 2014, funded by
650 the Templeton World Charity Foundation. The cost of open access publication was met by the Queen's University
651 Belfast.

652 **Author Contributions:** Just one author

653 **Conflicts of Interest:** The author declares no conflict of interest. The founding sponsors had no role in the design
 654 of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and in the
 655 decision to publish the results.

656 Abbreviations

657 The following abbreviations are used in this manuscript:

658	DFA	Derterminate Finite State Automaton
	DM	Determinate Machine
	FSA	Finite State Automaton
659	NDM	Non-Determinate Machine
	TM	Turing Machine
	UTM	Universal Turing Machine

660 Bibliography

- 661 1. Sharov, A.A. Functional Information: Towards Synthesis of Biosemiotics and Cybernetics. *Entropy* **2010**,
 662 12, 1050–1070.
- 663 2. Hazen, R.M. The emergence of patterning in life's origin and evolution. *Int. J. Develop. Biol.* **2009**,
 664 53, 683–692.
- 665 3. Farnsworth, K.; Nelson, J.; Gershenson, C. Living is Information Processing: From Molecules to Global
 666 Systems. *Acta Biotheor.* **2013**, 61, 203–222.
- 667 4. Weizsäcker, C.v. *Die Einheit der Natur [The unity of nature]*; Deutscher Taschenbuch Verlag: Munich,
 668 Germany, 1974.
- 669 5. Stonier, T. Information as a basic property of the universe. *Bio Systems* **1996**, 38, 135–140.
- 670 6. Devlin, K.J. *Logic and information*; Cambridge University Press: Cambridge, U.K., 1992.
- 671 7. Floridi, L. Information. In *The Blackwell Guide to the Philosophy of Computing and Information*; Floridi, L., Ed.;
 672 Blackwell Publishing Ltd, 2003; pp. 40–61.
- 673 8. Szostak, J.W. Functional information: Molecular messages. *Nature* **2003**, 423, 689–689.
- 674 9. Hazen, R.M.; Griffin, P.L.; Carothers, J.M.; Szostak, J.W. Functional information and the emergence of
 675 biocomplexity. *Proc. Natl. Acad. Sci. U. S. A.* **2007**, 104, 8574–8581.
- 676 10. Hagiya, M.; Aubert-Kato, N.; Wang, S.; Kobayashi, S. Molecular computers for molecular robots as hybrid
 677 systems. *Theor. Comp. Sci.* **2016**, 632, 4–20. Verification of Engineered Molecular Devices and Programs.
- 678 11. Kane, R. *A contemporary Introduction to free will.*; Oxford University Press.: Oxford, UK, 2005.
- 679 12. Westen, P. Getting the Fly out of the Bottle: The False Problem of Free Will and Determinism. *Buffalo*
 680 *Criminal Law Review.* **2005**, 8, 599–652.
- 681 13. Van Inwagen, P. *An Essay on Free Will.*; Oxford University Press.: Oxford, UK, 1983.
- 682 14. List, C. Free will, determinism, and the possibility of doing otherwise. *Noûs* **2014**, 48, 156–178.
- 683 15. Van Inwagen, P. Some Thoughts on An Essay on Free Will. *Harvard Rev. Phil.* **2015**, 22, 16–30.
- 684 16. Strawson, G. *Freedom and Belief*; Oxford University Press: Oxford, UK, 1986.
- 685 17. Kauffman, S.A. Autocatalytic sets of proteins. *J. Theor. Biol.* **1986**, 119, 1–24.
- 686 18. Kauffman, S.; Clayton, P. On emergence, agency, and organization. *Biol. Philos.* **2006**, 21, 501–521.
- 687 19. Ginsborg, H., *A Companion to Kant.*; Wiley-Blackwell, 2006; chapter Kant's biological teleology and its
 688 philosophical significance.
- 689 20. Hordijk, W.; Steel, M. Detecting autocatalytic, self-sustaining sets in chemical reaction systems. *J. Theor.*
 690 *Biol.* **2004**, 227, 451–461.
- 691 21. Hordijk, W.; Hein, J.; Steel, M. Autocatalytic Sets and the Origin of Life. *Entropy* **2010**, 12, 1733–1742.
- 692 22. Segré, D.; Ben-Eli, D.; Lancet, D. Compositional genomes: Prebiotic information transfer in mutually
 693 catalytic noncovalent assemblies. *Proc. Natl. Acad. Sci. U. S. A.* **2000**, 97, 4112–4117.
- 694 23. Luisi, P. Autopoiesis: a review and a reappraisal. *Naturwissenschaften* **2003**, 90, 49–59.
- 695 24. Davies, J.A. *Life Unfolding: how the human body creates itself.*; Oxford University Press: Oxford, UK, 2014.
- 696 25. Von Neumann, J.; Burks, A. *Theory of self-reproducing automata*; Illinois University Press: Chicago, IL.,USA.,
 697 1966.

- 698 26. Varela, F.; Maturana, H.; Uribe, R. Autopoiesis: the organization of living systems, its characterization and
699 a model. *Curr. Mod. Biol.* **1974**, *5*, 187–96.
- 700 27. Auletta, G.; Ellis, G.; Jaeger, L. Top-down causation by information control: from a philosophical
701 problem to a scientific research programme. *J. R. Soc. Interface* **2008**, *5*, 1159–1172,
702 [<http://rsif.royalsocietypublishing.org/content/5/27/1159.full.pdf>].
- 703 28. Jaeger, L.; Calkins, E.R. Downward causation by information control in micro-organisms. *Interface Focus*
704 **2012**, *2*, 26–41.
- 705 29. Cummins, R. Functional Analysis. *J. Philos.* **1975**, *72*, 741–765.
- 706 30. Farnsworth, K.D.; Albantakis, L.; Caruso, T. Unifying concepts of biological function from molecules to
707 ecosystems. *Oikos* - In review.
- 708 31. Lorenz, D.; Jeng, A.; Deem, M. The emergence of modularity in biological systems. *Phys. Life Rev.* **2011**,
709 *8*, 129–160.
- 710 32. Butterfield, J. Laws, causation and dynamics at different levels. *Interface Focus* **2012**, *2*, 101–114.
- 711 33. List, C. Levels: descriptive, explanatory, and ontological.
- 712 34. Ellis, G. Top-down causation and emergence: some comments on mechanisms. *Interface Focus* **2012**,
713 *2*, 126–140.
- 714 35. Ellis, G. On the nature of causation in complex systems. *Trans. R. Soc. S.Afr.* **2008**, *63*, 1–16.
- 715 36. Farnsworth, K.D.; Ellis, G.F.; Jaeger, L., *Living through Downward Causation*; Cambridge University Press,
716 2017; chapter 13, pp. 303–333.
- 717 37. Gilbert, S. *Developmental Biology*; Sinauer Associates: Sunderland, MA, USA., 2013.
- 718 38. Seeley, T. *Honeybee Democracy*; Princeton University Press: Princeton, USA., 2010.
- 719 39. Lineweaver, C.H.; Egan, C. Life, gravity and the second law of thermodynamics. *Phys. Life Rev.* **2008**,
720 *5*, 225–242.
- 721 40. Adami, C.; Ofria, C.; Collier, T. Evolution of biological complexity. *Proc. Natl. Acad. Sci. U. S. A.* **2000**,
722 *97*, 4463–4468.
- 723 41. Hoel, E.; Albantakis, L.; Marshall, W. Can the macro beat the micro? Integrated information across
724 spatiotemporal scales. *J. Conscious. Sci.* **2016**, *1*, niw012.
- 725 42. Marshall, W.; Albantakis, L.; Tononi, G. Black-boxing and cause-effect power. *ArXiv e-prints* **2016**,
726 [[arXiv:q-bio.NC/1608.03461](https://arxiv.org/abs/1608.03461)].
- 727 43. Mayr, E. Teleological and Teleonomic: A New Analysis. *Boston Studies Phil. Sci.* **1974**, pp. 91–117.
- 728 44. Walker, S.I. Top-down causation and the rise of information in the emergence of life. *Information* **2014**,
729 *5*, 424–439.
- 730 45. Wilson, D.; Wilson, E.O. Rethinking the Theoretical Foundation of Sociobiology. *Q. Rev. Biol.* **2007**,
731 *82*, 327–348, [<http://dx.doi.org/10.1086/522809>]. PMID: 18217526.
- 732 46. Wilson, E.O. *Sociobiology: The New Synthesis.*; Harvard University Press: Cambridge (MA), USA., 1975.
- 733 47. Danchin, A. Bacteria as computers making computers. *FEMS Microbiol. Rev.* **2009**, *33*, 3–26.
- 734 48. Walker, S.; Davies, P. The algorithmic origins of life. *J. R. Soc. Interface* **2013**, *10*.
- 735 49. Kauffman, S.A. *Investigations*; Oxford University Press, 2000.
- 736 50. Ptashne, M. Principles of a switch. *Nat. Chem. Biol.* **2011**, *7*, 484–487.
- 737 51. Moore, E.F., *Automata Studies*; Princeton University Press., 1956; chapter Gedanken-Experiments on
738 sequential machines., pp. 129–153.
- 739 52. Minsky, M.L. *Computation – Finite and Infinite Machines.*; Prentice Hall: Englewood Cliffs, NJ, U.S.A., 1967.
- 740 53. Zhang, J. Adaptive learning via selectionism and Bayesianism, Part 1: Connection between the two. *Neural*
741 *Networks.* **2009**, *22*, 220228.
- 742 54. Bruce, L.L.; Neary, T.J. The Limbic System of Tetrapods: A Comparative Analysis of Cortical and
743 Amygdalar Populations. *Brain Behav. Evol.* **1995**, *46*, 224–234.

