



ACADEMIC
PRESS

Consciousness and Cognition xxx (2002) xxx-xxx

**Consciousness
and
Cognition**

www.academicpress.com

Counterfactuals cannot count: A rejoinder to David Chalmers[☆]

Mark Bishop*

*Cybernetics Intelligence Research Group, Department of Cybernetics, University of Reading,
Berkshire, UK*

Received 3 July 2002

Abstract

9 The initial argument presented herein is not significantly original—it is a simple
10 reflection upon a notion of computation originally developed by Putnam (Putnam,
11 1988). See also Searle (1990) and criticised by Chalmers et al. (see Chalmers (1994,
12 1996a,b) see also the special issue, *What is Computation?*, in *Minds and machines*,
13 4:4, November 1994). In what follows, instead of seeking to justify Putnam's con-
14 clusion that every open system implements every Finite State Automaton (FSA) and
15 hence that psychological states of the brain cannot be functional states of a com-
16 puter, I will establish the weaker result that, over a finite time window every open
17 system implements the trace of FSA **Q**, as it executes program (**P**) on input (**I**). If
18 correct the resulting bold philosophical claim is that phenomenal states—such as
19 feelings and visual experiences—can never be understood or explained functionally.
20 © 2002 Published by Elsevier Science (USA).

21

[☆]The argument developed herein followed discussions with David Chalmers regarding the authors papers, 'Dancing with Pixies,' presented at ASSC4 (Brussels 2000) and 'Counterfactuals Can't Count,' presented at ASSC5 (North Carolina 2001). An earlier treatment of this material, including background to the debate, will appear as, Bishop, J. M. (2002). *Dancing with Pixies*. In J. Preston & J. M. Bishop (Eds.), *Views into the Chinese Room*, Oxford: OUP.

*Fax: +44-118-931-8220.

E-mail address: j.m.bishop@reading.ac.uk.

22 **1. Introduction**

23 In the first part of this paper I introduce Discrete State Machines (DSMs) and
24 show how, with input to them defined, their behaviour is described by a simple
25 unbranching sequence of state transitions analogous to that of an inputless FSA.
26 Then I review Putnam’s 1988 argument that purports to show how every open
27 physical system implements every inputless FSA. This argument is subsequently
28 applied to a robotic system that is claimed to instantiate genuine phenomenal states
29 as it operates—if the robot system instantiates such states then so must any open
30 system.

31 This type of argument has previously been criticised on the grounds that Putnam’s
32 open system does not correctly implement counterfactuals.¹ However, in a novel
33 extension I demonstrate that, in the execution of the robot’s control program with
34 known input, unless we allow non-executed state transition sequences to directly
35 influence the generation of phenomenal states, the lack of counterfactuals in Put-
36 nam’s open system cannot rescue functionalism; counterfactuals cannot count.

37 **2. Discrete state machines**

38 In his 1950 paper, *Computing Machinery and Intelligence*, Turing defined DSMs
39 as, “machines that move in sudden jumps or clicks from one quite definite state to
40 another” (Turing, 1950, p. 439), and explained that modern digital computers fall
41 within the class of them. An example DSM from Turing is that of a wheel machine
42 that clicks round through 120° once a second, but may be stopped by the application
43 of a lever-brake mechanism. In addition, if the machine stops in a pre-specified lo-
44 cation at one of the three possible positions, it will cause a lamp to come on. Input to
45 the machine is thus the position of the lever-brake, {brake on; brake off}, and the
46 output of the machine is the lamp state, {lamp on; lamp off}.

47 Such a machine can be described abstractly in the following manner. Its internal
48 (computational) states are labelled (arbitrarily) by a mapping function f that maps
49 from the physical state of the machine (i.e., what position the wheel is in) to the
50 computational state of the machine, $q \in \{q_1 q_2 q_3\}$. The input to the DSM, the brake
51 position, is described by an input signal, $i \in \{i_0 \text{ brake off}; i_1 \text{ brake on}\}$. Hence, the
52 next internal state of the machine is determined solely by its current state and its
53 current input as follows:

	q_1	q_2	q_3
i_0	q_2	q_3	q_1
i_1	q_1	q_2	q_3

55 With its output being determined by:

¹ Strong FSA conditional statements of the form, ‘if the FSA is in state a , it will transit into state b , however it finds itself in the first state.’

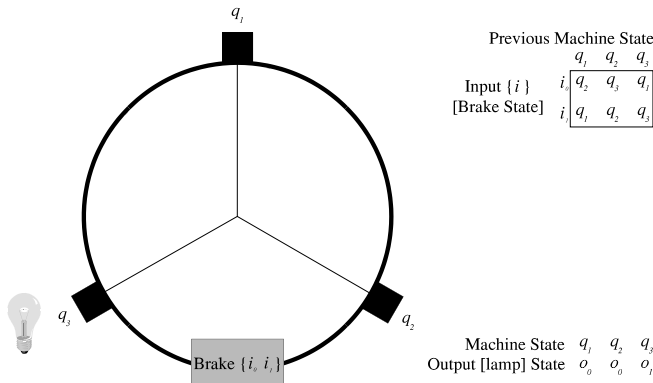


Fig. 1. Turing's, 'Discrete State Wheel Machine.'

State : q_1 q_2 q_3
 Output : o_0 o_0 o_1 [lamp on]

57 Thus, with input undefined, the state transition diagram in Fig. 1 describes the
 58 DSM's behaviour.

59 Fig. 2 has several branch points where the next state of the machine is determined
 60 by a state transition contingent on the current input. However, as shown in Fig. 3,
 61 for any specific input value there are no such branching state transitions. The ma-
 62 chine's output (lamp on/off) is determined purely by its initial state, ($q \in \{q_1 q_2 q_3\}$),
 63 and the system input value, ($i \in \{i_0 i_1\}$).

64 Knowledge of the specific input to the machine's state transition table (program)
 65 has thus collapsed its combinatorial structure. Further, over a given time period, say
 66 $[t_1 \dots t_7]$, all loops can be removed from the state diagram to form a linear path of
 67 state transits. The machine now functions, like clockwork, e.g.,

INPUT STATE 0: $\langle q_1 q_2 q_3 q_1 q_2 q_3 q_1 \rangle$ OR $\langle q_2 q_3 q_1 q_2 q_3 q_1 q_1 \rangle$
 OR $\langle q_3 q_1 q_2 q_3 q_1 q_2 q_3 \rangle$

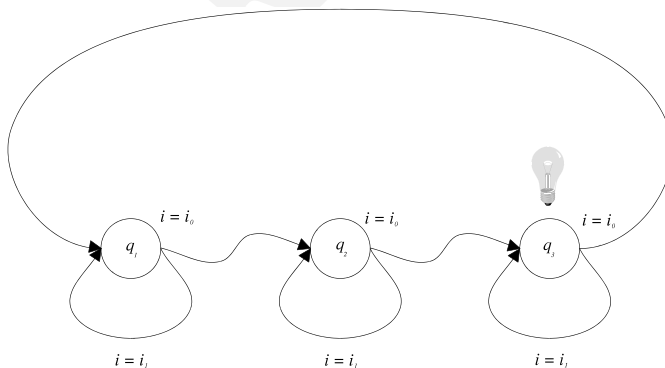


Fig. 2. State transition diagram of Turing's wheel machine—input undefined.

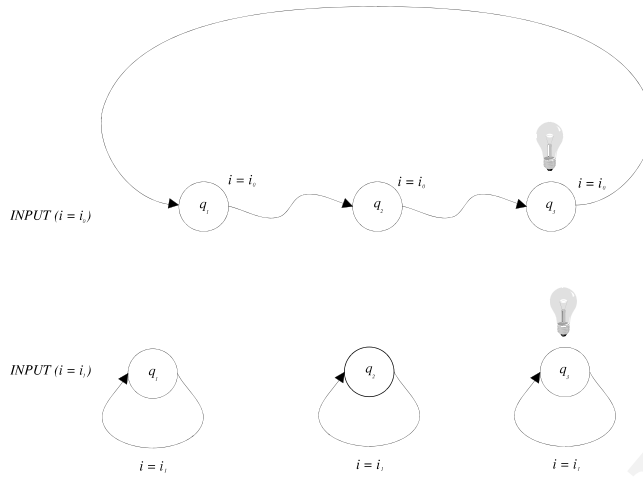


Fig. 3. State transition diagram of Turing's wheel machine—input defined.

69 **OR**

INPUT STATE 1: $\langle q_1 q_1 q_1 q_1 q_1 q_1 q_1 \rangle$ OR $\langle q_2 q_2 q_2 q_2 q_2 q_2 q_2 \rangle$
OR $\langle q_3 q_3 q_3 q_3 q_3 q_3 q_3 \rangle$

71 The following argument endeavours to show that for any Discrete State Machine
72 which, it is claimed, instantiates mental (phenomenal) states purely in virtue of its
73 execution of a suitable state transition table, *or program*, we can generate a corre-
74 sponding state transition sequence using any open physical system. That this conclu-
75 sion leads to a form of panpsychism, where all parts of matter involve
76 phenomenal consciousness, is clear, as if such state transition sequences are effec-
77 tively found in all open systems, then phenomenal states must be equally ubiquitous
78 in matter.

79 **3. Putnam's argument**

80 Tucked away in an appendix to Hilary Putnam's 1988 book, 'Representation and
81 Reality', is a short argument that endeavours to prove that any open physical system
82 is a realisation of every abstract inputless Finite State Automaton and hence that
83 Functionalism fails to provide an adequate foundation for the study of the mind.

84 Central to Putnam's original argument is the observation that every open physical
85 system, *S*, is in different maximal states² at every discrete instant and is characterised
86 by a discrete series of non-cyclic³ modal state transitions, $[s_1, s_2 \dots s_t \dots s_n]$. To
87 simplify the following discussion of Putnam's claim, I will replace Putnam's arbitrary

² Putnam (1988, p. 122).

³ Putnam (1988, p. 121).

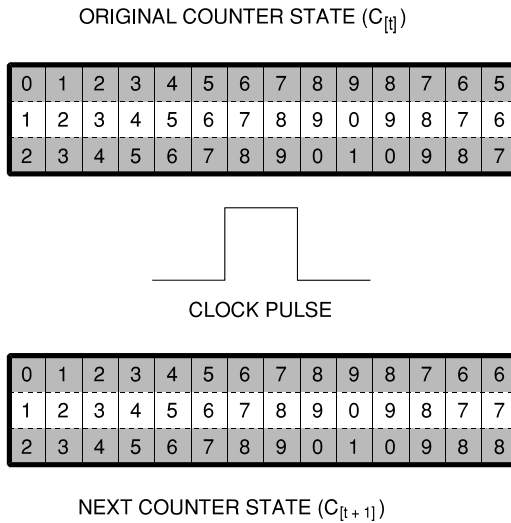


Fig. 4. A ‘non-cyclic’ counting machine.

88 open physical system, **S**, with an Infinite Counting Machine (ICM) generating the
 89 non-cyclic counter state sequence $[c_1, c_2 \dots c_t \dots c_n]$ in place of the physical system
 90 state sequence $[s_1, s_2 \dots s_t \dots s_n]$ ⁴ (see Fig. 4).

91 It is clear that given counter state $[c_k]$ at time $[t_k]$, it is trivial to predict its next
 92 state $[c_{k+1}]$ at time $[t_{k+1}]$ NB. This transition from state $[c_k]$ to $[c_{k+1}]$ is both regular
 93 and carries full modal force—that the counting machine is in state $[c_k]$ defines and
 94 contains the provision to force it to transit to $[c_{k+1}]$ at the next clock interval.

95 Any inputless FSA is characterised by its state transition table, defining, given its
 96 current state, its subsequent state. Imagine, without loss of generality, that the state
 97 transition table for FSA **Q** calls for the automaton to go through the following
 98 sequence of states in the interval $[t_1 \dots t_6]$:

$$\langle q_1 \ q_2 \ q_3 \ q_1 \ q_2 \ q_3 \rangle$$

100 Next let us suppose we are given a counting machine, **C**, which goes through the
 101 sequence of states, $[c_1, c_2, c_3, c_4, c_5, c_6]$ in the interval $[t_1 \dots t_6]$. We wish to find a
 102 mapping between counter states and FSA states such that, during the time interval
 103 under observation, the counting machine obeys **Q**’s state transition table by going

⁴ Putnam uses an open physical system simply as a device to generate a non-cyclic series of discrete states. Chalmers (1996a) argues that for such a system to reliably transit a sequence of states it must include a natural clock (such as a source of radioactive decay), however he also concedes that, “probably most physical systems satisfy such a requirement” (Chalmers, 1996a, p. 316). Nonetheless, to simplify this exposition of Putnam’s mapping I initially employ the conceptually simpler *ICM*, (an example of a finite counting machine being the milometer of a car). However, to fully generalise the results of the argument, in subsequent sections I will resort to using Putnam’s *open physical system*.

104 through a sequence of states, which the state mapping function will label
105 $[q_1 q_2 q_3 q_1 q_2 q_3]$.

106 4. Putnam's mapping

107 It is trivial to observe that if we map FSA state $[q_1]$ to the disjunction of counting
108 machine states, $[c_1 v c_4]$, FSA state $[q_2]$ to the disjunction of counting machine states,
109 $[c_2 v c_5]$ and FSA state $[q_3]$ to the disjunction of counting machine states, $[c_3 v c_6]$,
110 then the counting machine will fully implement \mathbf{Q} . Further, given any counting
111 machine state $[c_a] \in \{c_1, c_4\}$, at time $[t_a]$, we can predict the FSA will enter state $[q_2]$
112 at time $[t_b]$.

113 To show that being in state $[q_1]$ at time $[t_1]$ caused the counting machine to enter
114 state $[q_2]$ at $[t_2]$ we observe that at $[t_1]$ the counting machine is in state $[c_1]$, (which the
115 mapping function labels FSA state $[q_1]$), and that being in state $[c_1]$ at $[t_1]$ causes the
116 counting machine to enter state $[c_2]$, (which the mapping function labels FSA state
117 $[q_2]$) at $[t_2]$. Hence, given the current state of the counting machine at time $[t]$, we can
118 predict its future state and hence how the states of \mathbf{Q} evolve over the time interval
119 under observation.

120 Note, after Chalmers, that the counting system above will only implement a
121 particular execution run of the FSA. Hence, Chalmers remains unfazed at this result
122 because he states that inputless FSA's are simply an inappropriate formalism for a
123 Computationalist Theory of Mind.

124 To see the triviality, note that the state space of an inputless FSA will consist of a single unbran-
125 ching sequence of states ending in a cycle, or at best in a finite number of such sequences. The
126 latter possibility arises if there is no state from which every state is reachable. It is possible that
127 the various sequences will join at some point, but this is as far as the 'structure' of the state space
128 goes. This is a completely uninteresting kind of structure, as indeed is witnessed by the fact that
129 it is satisfied by a simple combination of a clock and dial. (Chalmers, 1996a, p. 318)

130 But Putnam extends his result to the case of FSA's with input and output, by
131 arguing that an FSA with input and output is realised by every open physical system
132 with the right input/output dependencies—if the physical system has the right input/
133 output then it instantiates the FSA correctly. Patently, this is a restriction on his
134 original claim, but nonetheless it remains a significant result that, if correct, suggests
135 functionalism implies behaviourism.⁵

136 Putnam's original argument, using an open physical system to generate a series of
137 non-repeating system states, (equivalent to the non-cyclic counting machine states
138 used above), runs as follows. For any arbitrary FSA, take an open physical system
139 with the right input/output dependencies. For example a rock with a number of
140 marks on it encoding the input vector, (\mathbf{I}) , (where (\mathbf{I}) encodes the finite set of input

⁵ "Any open system with the correct input/output dependencies implements the FSA with input/output. Hence, every FSA with input (\mathbf{I}) and output (\mathbf{O}) is implemented by any physical system with the same input/output dependencies. Hence mentality is contingent only on input and output and functionalism implies behaviourism." (Putnam, 1988, pp. 124–125).

141 values, $\{i_1, i_2 \dots i_n\}$) and another set of marks encoding the output vector (**O**), (where
142 (**O**) encodes the finite set of output values, $\{o_1, o_2 \dots o_n\}$). Associate rock state $[s_1]$ at
143 $[t_1]$ with the relevant initial state of the FSA, and rock state $[s_2]$ at $[t_2]$ with the
144 subsequent state of the FSA etc. Putnam claims it is clear that by mapping each FSA
145 state with the disjunction of associated rock states we ensure the system goes through
146 the relevant state sequence $[s_1, s_2, s_3]$ that corresponds, using this mapping, to the
147 relevant FSA state sequence, $[q_a, q_b, q_c]$, with system output encoded by the other
148 marks on the rock.

149 However, as for Turing's DSM, the addition of input now makes the formalism non-
150 trivial. There can now be branching in the execution trace, as the next FSA state is
151 contingent upon its current state and the system input. This gives the system a com-
152 binatorial structure. However, as Chalmers states, Putnam's revised construction does
153 not fully encapsulate this structure—rather it merely manifests one trace of the FSA
154 with a specific input/output dependency. However, as we observed with Turing's simple
155 three state DSM, when input to the system is defined over a specific time interval the
156 combinatorial state structure of the FSA collapses to a bounded linear path, which can
157 be generated using a Putnam style mapping and any open physical system.

158 5. Simulating phenomenal states over a finite time window

159 Consider a mobile robot whose behaviour is controlled by a program (**P**), acting
160 upon input (**I**), generating output (**O**), running on computer hardware (**Q**), over the
161 specified time interval $[t_1 \dots t_k]$. Assume that after switch-on at time $[t_1]$, the robot
162 experiences a series of phenomenal states until it is switched off at time $[t_k]$.

163 During the specified interval, $[t_1 \dots t_k]$, the robot's input states are defined by data
164 from its sensors forming the input set (**I**) = $\{i_1 i_2 i_3 \dots i_k\}$ with its output states
165 defining the actuator commands controlling its external behaviour, which together
166 form the output set (**O**) = $\{o_1 o_2 o_3 \dots o_k\}$. Let us now concede that the robot's
167 control program, (**P**), instantiates a series of phenomenal states in **Q**, caused by
168 interaction with its environment as the program executes. What is it in the robot that
169 manifests this property? Unless we allow these phenomenal states to extend outside
170 of the robots computational hardware, **Q**, they must be manifested solely by it. That
171 is, the claimed phenomenal properties of the system must be solely realised by **Q**
172 (**P,I**); the computer **Q**, executing program (**P**), on input (**I**).

173 Now as **Q** executes (**P**) over the interval $[t_1 \dots t_k]$, **Q** (**P, I**) generates a specific set
174 of computational states, **S**, (**S** = $\{s_1, s_2 \dots s_r\}$), at the discrete clock intervals of the
175 CPU, **Q**. Due to the universal realisability of Turing Machine programs,⁶ the
176 particular underlying computational engine, **Q**, is irrelevant to the generation of
177 the computational states $\{s_1, s_2 \dots s_r\}$. Whether **Q** had at its heart a CPU made
178 from toilet rolls⁷ or the latest Intel processor, the associated computational states

⁶ All digital computers are in a sense equivalent—the underlying hardware that executes a specific computation does not affect the result of that computation.

⁷ Weizenbaum (1976, p. 51).

179 are the same and it is **Q**'s generation of these states that results in its instantiation
180 of phenomenal experience.

181 However, we have already seen from Putnam how the computational states re-
182 sulting from the execution of any given FSA, **Q**, (with specific input/output) can be
183 mapped onto an open physical system, (e.g., a simple counting machine). Thus, by
184 relaxing the requirement that the physical system instantiates the full combinatorial
185 state structure of a program with general input, to the relatively trivial requirement
186 that it instantiate the correct state transitions for a specific execution trace, we
187 sidestep the need for an exponential increase in state space. Thus, with the input set
188 defined, over the specified time interval $[t_1 \dots t_k]$ we are able to fully generate the
189 computational states governing the robot's behaviour and hence instantiate its
190 phenomenal states, via a simple counting machine or, following Putnam, any open
191 physical system.

192 6. Counterfactuals can't count

193 The above argument has been criticised by David Chalmers and others⁸ on the
194 grounds that the Putnam FSA analogue, (the open physical system/non-cyclic
195 counter described earlier), does not genuinely replicate an FSA with input. In par-
196 ticular, it lacks the ability to correctly implement counterfactuals. Specifically, 'if a
197 state-transition in the machine table does not play a role in that run, then it will
198 correspond to nothing in the structure of the physical system. The system therefore
199 fails to implement the FSA.' It would seem that even though specific computational
200 states may not be entered on a specific execution run of the FSA, the mere possibility
201 that they could be if system input was different is required in a genuine functional
202 isomorph of the FSA.

203 However, in the next section, by employing a variation on Chalmers' Fading
204 Qualia Argument (FQA), I will show that in the context of an FSA with input fixed a
205 priori, state transition behaviour reduces to that of an inputless FSA, and Putnam's
206 FSA open system analogue fully holds.

207 7. A clockwork robot

208 Consider the operation of two robots over the time interval $[t_1 \dots t_k]$, with defined
209 input (**I**), ($\mathbf{I} = [i_1 \dots i_k]$). A program designed to Chalmers' specification, replicating
210 the fine-grained functional organisation of a system known to have phenomenal
211 states, controls one robot, R_1 ; the second, R_2 , generates the computational states of
212 R_1 via an open physical system using Putnam's mapping.

213 Although the input to both robots is fully defined a priori, ($\mathbf{I} = [i_1 \dots i_k]$), and
214 hence the series of computational states transited by both systems over the specified

⁸ In discussions with: Peter Fletcher at the University of Keele, (UK, 1999); David Chalmers at ASSC5 (North Carolina 2001).

215 time interval is identical, the *counterfactual hypothesis* is that only R_1 would expe-
216 rience phenomenal states, as only R_1 correctly implements counterfactuals and
217 maintains behaviour that is sensitive to the input as its program execution trace,⁹ T_{R_1}
218 (\mathbf{I}), evolves.

219 However, consider what happens if, contingent upon the input (\mathbf{I}), at each con-
220 ditional branch point in the execution trace $T_{R_1}(\mathbf{I})$, we delete a state transition se-
221 quence that is not entered,¹⁰ then iteratively repeat this procedure until Chalmers'
222 robot, R_1 , with full input-sensitivity, is step-by-step transformed into a copy of the
223 second robot R_2 , such that its resulting behaviour is also determined solely by a series
224 of linear state transitions. We can imagine that throughout this replacement pro-
225 cedure the robot repeatedly reports the colour of a red square fixed within its sensor
226 field; that is, the robot input vector (\mathbf{I}) remains constant and encodes the red square,
227 ($\mathbf{I} = [\text{red square}]$).

228 Between R_1 and R_2 there are a number of intermediate robots
229 [$R_1 \dots R_{(i-1)} R_{(i)} R_{(i+1)} \dots R_2$]—what is it like to be them? As we transform R_1 into R_2 , how
230 does its phenomenal perception vary? Imagine that initially R_1 was having bright ex-
231 perience of the red square, the counterfactual hypothesis would be that at some point in
232 the transformation $R_{(i)}$'s phenomenal experience must cease to be bright red. However,
233 the only difference between $R_{(i-1)}$ and $R_{(i)}$, given the bright red square remains fixed as
234 input, is that a sequence of states that was not transited in $T_{R_{(i-1)}}(\mathbf{I})$, has been deleted
235 from $R_{(i)}$.

236 It is clear that this scenario is implausible, for otherwise we have a robot, $R_{(i)}$,
237 whose phenomenal experience of input (\mathbf{I}) is in some way contingent upon the mere
238 presence of potential state transition sequences that are never transited during FSA
239 operation.¹¹ Hence the counterfactual hypothesis is false and the phenomenal ex-
240 perience of R_1 and R_2 must be the same.

241 Yet is this version of the FQA valid? David Chalmers has argued that it is not. In
242 contrast to the original FQA, in this scenario although the first robot, R_1 , is sensitive to
243 its input, the second, R_2 , is clearly not ('it merely acts like a clockwork toy').

244 It is clear that removing sensitivity to input per se cannot affect $R_{(i)}$'s phenomenal
245 experience, for consider what would happen to R_1 if the link between its frame store¹²
246 and its visual sensor is damaged, such that its frame store constantly maintains digital
247 representation of a red square, irrespective of the true image that falls on the optical

⁹ Program execution trace, $T_{R_1}(\mathbf{I})$, is the sequence of state transitions generated by the FSA controlling robot R_1 , given input (\mathbf{I}).

¹⁰ This is achieved by replacing one input-sensitive branching state transition, (cf. Fig. 2), with a simple linear state transition contingent on the input, (cf. Fig. 3).

¹¹ In other words, the mere presence of segments of program code that are not executed by the robot's computer would have direct casual influence over the phenomenal states instantiated in the machine.

¹² A device that maintains a digital representation of an image obtained by a visual sensor such as a TV camera.

248 sensor. This would result in R_1 , like R_2 , becoming insensitive to its visual input.¹³ When
249 asked, R_1 will now ‘act like clockwork’ and always report that the object in its visual
250 field is a red square, irrespective of its true shape and hue.

251 The lack of input-sensitivity will either deflate the phenomenal experience of R_1 or
252 have no effect (with any phenomenal states in the latter case analogous to the human
253 experience of a red hallucination). But as R_1 ’s control program is unchanged and the
254 data it reads from its frame store is exactly the same as usual, (a set of binary numbers),
255 R_1 will continue to function as normal and its phenomenal experience must be un-
256 changed. Hence, in the execution of an FSA with fixed constant input, the lack of input-
257 sensitivity per se will not influence the phenomenal states instantiated by the system.

258 8. Conclusion

259 For any computing machine, \mathbf{Q} , executing program (\mathbf{P}), with known input (\mathbf{I}),
260 over the specified time interval $[t_1 \dots t_k]$, only a formal [and repeatable] series of state
261 transitions occur within its hardware. The generation of the state transitions must be
262 responsible for the generation of the machines phenomenal properties. In this paper
263 we have seen why, following Turing’s observations on universal realisability, the
264 underlying hardware that instantiates computational state transitions is unimportant
265 and hence, after Putnam, that in principle a series of such transitions could be im-
266 plemented by any open physical system. Finally we observed that, with input to \mathbf{Q}
267 (\mathbf{P} , \mathbf{I}) defined, the mere potential for counterfactual reasoning,¹⁴ cannot influence the
268 phenomenal states instantiated by the system. Thus, over a finite interval, if a
269 computational system \mathbf{Q} (\mathbf{P} , \mathbf{I}) does instantiate genuine phenomenal states, then so
270 will any open physical system; implying functionalism will fail, crediting too many
271 things with the phenomenal quality of mind.

272 Acknowledgments

273 The author gratefully thanks Dr. Slawomir Nasuto and the anonymous reviewers
274 for their helpful comments on the original manuscript.

275 References

- 276 Chalmers, D. (1994). On implementing a computation. *Minds and Machines*, 4, 391–402.
277 Chalmers, D. (1996a). Does a rock implement every finite-state automaton. *Synthese*, 108,
278 309–333.

¹³ In the earlier discussion R_1 is incrementally modified but its input data is constant. Here the FSA functions as designed, but its incoming data is modified, (it is forced to take on a specified value irrespective of the image on the robots optical sensor).

¹⁴ Strong FSA conditional state transitions.

- 279 Chalmers, D. (1996b). *The conscious mind: in search of a fundamental theory*. Oxford: Oxford
280 University.
- 281 Putnam, H. (1988). *Representation and reality*. Bradford Books.
- 282 Searle, J. (1990). Is the brain a digital computer. *Proceedings of the American Philosophical*
283 *Association*, 64, 21–37.
- 284 Turing, A. (1950). Computing machinery and intelligence. *MIND*, 49, 433–460.
- 285 Weizenbaum, J. (1976). *Computer power and human reason: from judgement to calculation*. San
286 Francisco: W.H. Freeman.

UNCORRECTED PROOF