
Multiple realizations of the mental states: hunting for plausible chimeras

VINCENZO G. FIORE

1 The framework of multiple realizability

The key elements characterising the functionalist approach to mind studies are commonly identified (e.g. see [5]) with claims concerning:

1. The cognitive creatures' essential feature (they are all computational systems);
2. The object of the research in the fields of cognitive psychology and artificial intelligence (abstract functional states and novel physical realizations for these states respectively);
3. The irreducibility and consequently the autonomy of special sciences;
4. The inefficacy of the empirical research on the neural structure, because of the merely contingent relation established between the neural structure and the functional states it realizes.

The objective of this paper is to support a reductionist perspective in mind studies, disputing the soundness of the claims 3 and 4 in particular. Therefore, since it is easy to concede that the theory of multiple realizability of mental states plays the role of the hub, binding all the four claims one another, this paper aims at showing the weaknesses of the grounds on which the theory has been built.

The Multiple Realizability Theory (MRT) has been first formalized in the late sixties by Hilary Putnam in a famous series of papers (for a collection see [12]). In the article commonly recognised as the most representative of that period [11], it is assumed that every animal, independently of the species it belongs to, is capable of feeling pain: the mental state of pain is not species-specific. Therefore the identification of the mental state with a certain C-Fiber activation (or any other neural correlate) leads to the conclusion that all species should be found sharing the same neural structure and the same neural activation at the right moment. Even if we consider that parallel

28 evolution might lead to the same physical structure, once the argument is
 29 extended to other psychological predicates (such as, for instance, hunger or
 30 sexual attraction), it becomes *overwhelmingly plausible* (Putnam's words)
 31 that these multiple realizations across species simply cannot be explained
 32 in terms of a theory grounded on the identity between mental and physical
 33 states. After all, even if parallel evolution could be proved in all known
 34 creatures, the conceivability of artificial silicon based systems capable of
 35 feeling pain, would definitively discard any attempt to establish an identity.
 36 Putnam's famous proposal is then to conceive a different approach to the
 37 mind, grounding it on the concept of a virtual machine analogous to the
 38 Turing Machine, but characterised by a few strategic differences.

It is useful to remind briefly what these devices are: a Turing machine (TM) is a computational -serial- device that is instructed by a program (set of instructions) to process a symbolic input in order to give a symbolic output as a result. These processes may have the following schematic representation:

$$\{x_1, x_2, x_3, \dots x_n\} \rightarrow A \rightarrow B \rightarrow C \rightarrow D \rightarrow \dots \rightarrow [final\ status]$$

39 The input assigns a value to each of the n variables $\{x_1, x_2, x_3, \dots x_n\}$, then
 40 the virtual machine computes these values as it is described by its set of
 41 instructions, reaching its first state (A). The new state gives life to a new
 42 series of processes that allows the machine to change again state in favour
 43 of the second one (B): the operation is replicated until the virtual machine
 44 reaches the final state described by the instructions in relation to the values
 45 assigned to the variables.

46 This mechanism implies that a TM is characterised by an assignment of
 47 probabilities 1 or 0 to every transition. On the contrary, if the instructions
 48 allow the machine to change its status from the original one to a series
 49 of target ones, with probabilities assigned to each of them, (e.g. starting
 50 from the functional state A the machine may change in favour of B with
 51 30% of chances or C with 70%) then the machine is called Probabilistic
 52 Automaton. Finally, there are devices capable of processing sets of inputs
 53 in order to generate new sets of instructions: this ability allows simulating
 54 any possible TM generating a so-called Universal Turing Machine (UTM). In
 55 other words, the UTM is directly programmed by the input, which instructs
 56 the machine about the processes to apply thenceforth. The MRT assumes
 57 that the combination between a probabilistic automaton and a UTM gives
 58 in return a virtual device whose processes are consistent with the living
 59 beings' ones.

60 All these devices (TM, UTM and probabilistic automaton) are known
 61 as virtual machines because of their nature which makes them completely

62 independent of any specific physical structure: it doesn't matter if the com-
 63 putation required by the set of instructions is performed by a neural system,
 64 a CPU or a series of cogs wheels. The focus is on the functional organization
 65 realized by the device (i.e. the instructions concerning its state transitions)
 66 and the functional state it can consequently reach, once the device has re-
 67 ceived a specific symbolic input. Furthermore, since the states are also
 68 independent, it is not even necessary for two systems to be functionally
 69 isomorphic (i.e. it is not necessary that they realize the same set of instruc-
 70 tions) to reach the same state: different programs may lead to the same
 71 functional state.

72 In conclusion, the MRT entails that two generic neural structures A and
 73 B may realize a mental state M, but they can never be identified with
 74 the mental state itself: the relation between the physical system and its
 75 mental realizations is always contingent and there can be infinite physically
 76 different systems realizing the same mental state. The focus changes from
 77 the reductionist study of the neural correlate to the functionalist study of
 78 the realized functions¹.

79 Putnam's early argument has been originally applied to different neural
 80 structures belonging to different species, but few years later Jerry Fodor
 81 [7, 8] generalised the value of the MRT, presenting his assumption as the
 82 necessary consequence of Putnam's conclusions. The generalised version of
 83 the MRT has started appealing to the 70s studies on brain mapping and to
 84 the notions of neural degeneracy and plasticity: the key argument coming
 85 from these studies is that the nervous system of higher organisms is able to
 86 accomplish a single psychological task in a wide variety of ways by means
 87 of several neurological parts of the whole structure. As a consequence,
 88 the relation between physical and mental states proves to be contingent
 89 even when it is applied to the same species or a single neural system²:
 90 time becomes a legitimate variable to take into account when considering
 91 the contingency of the causal relation between the physical system (the
 92 implementer) and the functional state (the implemented).

93 **2 The computability issue and the overestimation of** 94 **the UTM**

95 The superimposition of the processes performed by a virtual machine on
 96 the ones realized by cognitive organisms has been attractive since the very

¹Subsequent articles (e.g. see [2] or i [12, §14]) have also dealt with the problem of the realization of more than a single functional state (or psychological predicate) at the same time. The solution proposed assumes complex living beings are able of realizing the processes of several virtual machines at the same time (i.e. in parallel).

²E.g. a single human being realizes the same mental state of pain during childhood and adulthood, despite the differences characterising the same neural structure in the two periods.

97 beginning: even those who have tried to discard the functionalist approach
98 have rarely questioned the argument of the multiple realizations of mental
99 states and have preferred to focus their attention on the implications the
100 theory has on reductionism [5, 9, 10, 4]. A few exceptions are represented
101 by those [17, 15, 1] who have challenged the likelihood of the argument
102 by means of theoretical reasoning or stressing the failures of the predic-
103 tions implied the generalised MRT. Nonetheless, I think a computational
104 approach to this matter has been surprisingly ignored: the theory relies on
105 the identification of the mind with the TM; should this identification be
106 computationally inadequate, the MRT would be proved ill-grounded. As a
107 matter of fact, there are three reasons that lead to this conclusion.

108 The first reason is the limited range of Turing-computable algorithms. To
109 put it simple, the computational capacities of a TM are widely overestimated
110 and they are usually erroneously attributed to Turing himself. There is a
111 huge list of philosophical misconceptions about Turing's virtual machine [6]
112 and they are all grounded on the erroneous assumption that in his articles
113 Turing may have mathematically demonstrated how a UTM can compute
114 any algorithm (i.e. the mathematical function that formally describes the
115 set of instructions or program of the virtual machine) performed by any
116 other machine with any architecture, given enough time and memory.

117 What Turing did demonstrate is that a UTM can realize any algorithm
118 characterised by the following requirements (which define the 'mechanical
119 method'):

- 120 1. finite number of exact instructions (each instruction expressed with
121 a finite number of symbols) to make the machine change from one
122 functional state to the following one.
- 123 2. Finite number of state transitions to produce the expected result.
- 124 3. In principle, a human being can carry it out only aided by paper and
125 pencil.
- 126 4. It does not require insight or ingenuity to be carried out³.

127 For the purpose of this article, it is sufficient to point out that the set of
128 hypothetic algorithms realized by any TM is countable, that is to say, it is
129 characterized by the same order of infinite of the integers. On the contrary,
130 the number of all the hypothetic computable algorithms is uncountable (i.e.
131 of a higher order of infinite): hence, there is an infinite number of algorithms

³These notions have a formal and rigorous equivalent[16, 3]: for the purpose of this paper it is sufficient to refer to their informal version.

132 which have a mathematical description and cannot be realized by a UTM,
133 even if they are realized by differently structured systems.

134 If the algorithms implemented by neural systems are not found to meet at
135 least one of the four requirements for Turing-computability, it must be con-
136 cluded that a UTM may not simulate or even describe information processes
137 in living beings. Consequently, it is necessary to study the way biological
138 neural systems process their data, before formulating any hypothesis about
139 the possibility to realize such processes by means of a virtual machine. Un-
140 der these circumstances, the hypothesis of multiple realizations of mental
141 processes may be empirically falsified: MRT cannot be established a priori.

142 It may be argued that even if we could find out that neural systems do
143 not realize Turing-computable algorithms, this finding by itself would not
144 be enough to discard multiple realizability. A new hypothetical and more
145 powerful virtual machine might be conceived: different from the known Tur-
146 ing machines, it might widen the range of realizable algorithms, overcoming
147 some of, if not all, the weak points of the classic machines.

148 Nonetheless, it seems that such a powerful virtual machine is unlikely to
149 come and it is usually considered mathematically implausible⁴. Even if it
150 were plausible, this objection would not lead far from the prospected path:
151 these new hypothetic systems would not be asked to simulate a generic
152 new set of algorithms but those specific of the parallel distributed -neural-
153 systems. Once again, in order to be sure that the proper set of algorithms
154 is part of the domain of these new machines (proving the soundness of
155 MRT), it would be necessary to know beforehand what sort of algorithms
156 are implemented by neural systems.

157 This conclusion leads to the second reasoning against the plausibility of
158 the MRT. There is a particular causal relation between the physical struc-
159 ture of a neural system and the algorithm it implements: a neural network
160 realizes a sheaf of sets of mathematical functions⁵ defined by its architecture
161 and by the computation performed by each single node of the network. The
162 values assigned to the other variables, such as the weights of the synapses

⁴The existence and the features of devices that may result to be able to implement such Turing-incomputable algorithms have been debated at least for five decades. An essential bibliography and a brief account of this debate can be found in section two of the cited Copeland's article [6]. As a matter of fact, the probabilistic automaton already represents a virtual machine which is able to realize a wider set of algorithms, if compared to a TM. I mainly refer to the TM for the convenience of the reasoning, but the criticism is valid for the probabilistic automaton as well: the set of algorithms realized is still countable and the algorithms themselves are characterized by similar features.

⁵E.g. the equation ($ax + by = k$) describes a sheaf of straight lines. If we fix the constants (in this case: a, b, k) attributing them a value, the result is the equation of a single straight line (e.g. $2x + 3y = 1$). A set of straight lines describes the equations combined in single or multiple systems.

163 (i.e. the electrochemical conductivity of the synapses), fix the constants
 164 for any specific set of algorithms within this sheaf. Every modification in
 165 the architecture of the network or in the processes of the single nodes leads
 166 to a system that can or cannot solve a specific given task⁶.

167 If we use simple connectionist models, the sheaf of algorithms imple-
 168 mented can be mathematically described with ease: in these conditions, the
 169 analysis of the relation between the neural structure and the implemented
 170 algorithm makes us conclude that the former has a causal influence on the
 171 latter. Nonetheless, even if the systems show a higher order of complex-
 172 ity (such as those proper of biological networks), it is possible to have an
 173 idea of the sheaf of algorithms determined by the architecture, especially
 174 considering that, though extremely complex, single neurons compute their
 175 electrochemical signals in a way that can be described by adequate mathe-
 176 matical functions. In a few words, different neural systems realize different
 177 algorithms, require different amount of energy and time to perform the
 178 same task and -due to differences in vector conversion- differ in the way the
 179 information is encoded or stored, in the categories developed and in their
 180 resistance to physical damages. Thus, mathematical analysis of neural sys-
 181 tems is telling us a different story from the one told by the MRT: in order
 182 to be able to process information -precisely- in the same way, two neural
 183 systems must be physically identical (i.e. two biological neural systems can
 184 hardly ever be functionally isomorphic due to the known structural differ-
 185 ences across species and within the same one).

186 It is still possible to claim that whether or not two neural systems may
 187 perfectly match their processes implementing the same algorithm, this would
 188 not affect the hypothesis that a serial device may be conceived realizing neu-
 189 ral processes. Once a probabilistic automaton were shown simulating the
 190 information processes of a neural system, the possibility to separate single
 191 states in the virtual machine would make it irrelevant for the MRT the
 192 whole second reasoning. Yet, the problem with this criticism is that it does
 193 not consider both the arguments so far described at the same time:

- 194 A. Whether or not a virtual machine may realize the set of instructions
 195 implemented by a neural system can only be established a posteriori.
- 196 B. The physical structure in neural systems is directly responsible for the
 197 processes implemented.

198 The two premises A and B lead inevitably to:

⁶The logical operator XOR is often cited in literature: it is known that there is no way to realize this computation with a single layer neural network (e.g. see [14, chap. 19, sect. 3].

199 C. In order to support an anti-reductionist path (MRT), it is necessary
200 to use a reductionist strategy, seeking the knowledge concerning the
201 processes realized by a neural system.

202 When everything is taken into consideration, the proof in favour of the
203 multiple realizability of the mental states would be reached after it had
204 become irrelevant.

205 The third reason against the plausibility of the MRT is grounded on
206 the computational inadequacy of serial systems in simulating the unique
207 features of biological neural systems. Biological systems deal with contin-
208 uous and infinite inputs, processes and outputs, processing information in
209 a flow; on the contrary, a virtual machine necessarily works with discrete
210 and finite data and state transitions, following a step-by-step procedure.
211 External data can reprogram a UTM to make it change its processes (once
212 the input has changed the set of instructions, the device can also apply its
213 rules to previously incomputable data), but the neural systems are able to
214 change their processes both depending on and independently of the input.
215 For instance, biological systems based on neural structures require a specific
216 amount of energies in order to activate their systems: a lack of energy mod-
217 ifies the computational processes by means of a change in the computation
218 performed in the single neurons of the network. This change takes place
219 independently of both the awareness and the perception of such a lack in
220 the organism. This feature is not limited to the energy requirements: any
221 physical alteration⁷ directly modifies the way the information is processed
222 by the system, but cannot be considered as part of the input.

223 A simulation with a Universal Turing machine can hardly give an account
224 of these phenomena, despite the fact that they are very frequent in all
225 living beings based on neural systems. Interestingly, Fodor [7] has used the
226 argument of plasticity and degeneracy to propose his generalised version
227 of the theory, but I think that this argument can be of use also against
228 the virtual machine hypothesis, at least until these systems will be able to
229 realize algorithms which can only be reprogrammed by input information.

230 Lastly, such differences make the parallel neural systems more robust
231 in respect of time and energy requirements: if the processes are suddenly
232 interrupted due to a lack of time, these systems are still able to give an
233 output, even if it will probably differ from the one the system would have
234 reached having sufficient amount of time. On the contrary, the mechanical
235 method implies that a serial system needs to follow all the given instructions

⁷E.g. structural damages or any other alteration of the neural architecture, chemical or electrical interference in electrochemical synapses, modification of the metabolic state of the neurons, etc.

236 in order to perform its transition among states: the lack of the time required
 237 to accomplish it would cause a failure in giving an output.

238 **3 Making it through the MRT**

239 It may be argued that it is here discussed the multiple realization of a
 240 whole set of instructions, but the object of the MRT is a single, indepen-
 241 dent and isolated functional state, which has its equivalent in the mental
 242 state/psychological predicate of a living being. Nonetheless, the supposed
 243 isolation of single psychological predicates such as pain, hunger, etc. is ac-
 244 ceptable within the context of the known virtual machines, such as the UTM
 245 and the probabilistic automaton: these machines are characterised by serial
 246 processes and therefore allow the existence of autonomous functional states.
 247 Once the identification of the mind with virtual machines is disputed, the
 248 existence of states of this sort in the mind is challenged too: our self-beliefs
 249 about them may be misleading.

250 Let us push this line of thought a little farther. This article has outlined
 251 the following proportion:

252 *Set of instruction: Turing machine = algorithm: system whose processes*
 253 *are mathematically describable*

254 It may be argued that this proportion implies the following:

255 *Functional state: Turing machine = assignation of values to all variables*
 256 *in the algorithm: system whose processes are mathematically describable*

257 In the set of parallel neural systems (which is a subset of the mathemat-
 258 ically describable systems), this proportion would imply that a particular
 259 kind of activation pattern would take the place of the third term in the
 260 second proportion. Though different from the 'C-nerve activation' correctly
 261 defined as *philosopher's fiction*[1], this would be anyway a completely theo-
 262 retical object: a sort of photography of the entire structure, taking into ac-
 263 count the whole network, the activation and metabolic status of all neurons
 264 and the disposition of every synapse to propagate its signals. Consequently,
 265 any change in any of the variables involved, would generate a different assign-
 266 ation to the variables as well as a different mental state, a conclusion that
 267 may seem to lead to an unusable theoretical object.

268 The problem is that biological neural networks are dynamical informa-
 269 tion processing systems, and consequently this perspective brings forth the
 270 concept of a theoretical object (the photography of the whole structure)
 271 characterised by an unavoidable incoherence. If the new definitions imply
 272 a concept of mental state which is both unusable and incoherent, then it
 273 seems it would be a good idea to discard the whole thesis, on the basis of
 274 its implications.

275 I think this is not a good reasoning: an analogy with the field of analysis

276 in mathematics should help in this case. A sheaf of straight lines can be
277 studied both independently of the assignations of values to its constants and
278 after the partial or complete assignation of the same values; the variables
279 also contribute to locate specific parts or single points on the line analysed.
280 As a consequence, it is perfectly plausible to imagine general rules that can
281 be applied to parallel neural systems (e.g. the computation performed by
282 a single neuron is almost the same in every organism showing a central or
283 distributed neural system: this is the assignation of value to a constant),
284 other rules that are species specific (the macro structure of the neural net-
285 work shows its similarities) and finally those rules which are single-structure
286 specific and vary within a single organism depending on its natural devel-
287 opment, experience and accidents. The use of the fine and coarse grain of
288 analysis [1], should make it possible to relate the new born theoretical men-
289 tal states — indeed a dynamic concept, far from the static serial equivalent,
290 but still usable- to the variances here described across species or within the
291 single organism.

292 This use of the mathematical descriptions does not lead to a hyper local
293 reductionism: the single events in the flow of continuous processes of the
294 system are still comparable within the same species with an acceptable fine
295 grain of analysis and the tool that allows such a comparison relies again
296 in the mathematical description of the algorithms realised by the neural
297 processes. Furthermore, there are many advantages in pursuing the use of
298 this tool to understand mind processes. The algorithms describe the way
299 every possible signal is computed by a system: they are not influenced by
300 the presence of a specific stimulus or a combination of stimuli, neither they
301 rely on the analysis of visible behaviours or other forms of output. As
302 it was originally conceived by Putnam concerning the set of instructions
303 of a probabilistic automaton, the specific study of the algorithms imple-
304 mented by neural system would allow to describe every possible process
305 these system perform in each of their layers, reaching important results in
306 the understanding of the observable and hidden phenomena⁸.

307 4 Conclusions

308 This paper states a methodological problem. There is no computational
309 device able to realize all the uncountable possible algorithms: as a conse-
310 quence, if the object of mind studies are the psychological predicates, it is
311 necessary to study the specific processes that generate them. Whether or
312 not these will result to be multiply realized, the computational study of

⁸Along this path, the main obstacle is represented by the epistemic indeterminacy due to the order of complexity of the biological neural systems, but I assume that grounding the models on the findings in neuroscience, a better explanatory value will be granted.

313 neural structures is the necessary first step of a realistic approach to the
314 mind. Furthermore, contrary to what expected by the MRT, the more sci-
315 ence gives us tools to investigate neural systems, the more it seems that the
316 processes they implement are supervened by the physical matter and are
317 characterised by a series of unique features.

318 Whenever the processes realized by a particular system are inaccessible,
319 the only way to attempt an analysis consists in assuming that another sys-
320 tem, whose processes are accessible, is realizing some of the processes of the
321 first inaccessible system. This procedure creates a useful analogy allowing
322 an analysis narrowed to a part of the whole set of processes of the acces-
323 sible system: as a consequence, the new aimed description is partial and
324 indirect, because it refers to the supposed analogous system rather than to
325 the original one.

326 My claim is that when multiple realizability is applied to neural systems,
327 it is useful to conceive it as a tool giving access to incomplete descriptions of
328 the psychological predicates: a similar constraint does not entail to discard
329 the procedure as a whole, because there are still cases in which there is no or
330 little access to complete descriptions. Nevertheless, if a complete description
331 is accessible or if a better analogy is established (due to an accessible system
332 which is closer to the unaccessible one), then the new description must be
333 preferred to the partial one formerly achieved. In the field of mind studies,
334 in the past few years, the mental processes are becoming more and more
335 accessible and consequently new descriptions will be formalized thanks to
336 this change: on this new ground, new explanatory theories will be built,
337 showing substantial divergence if compared with the ones formerly inferred
338 on the ground of the MRT.

339 In the attempt to save the MRT from Shapiro's remarks [15], Rosenberg
340 has stated that this theory has been proposed to explain *the absence of dis-*
341 *coverable psychophysical laws in a way compatible with physicalism*[13]. It
342 seems today that we are moving towards the finding of these laws: should
343 this happen by means of the mathematical description of the processes re-
344 alised by the neural systems, the prediction here supported is that the mul-
345 tiple realizability tool will see the fields it has been applied so far restrained,
346 in favour of the new tools.

347 BIBLIOGRAPHY

- 348 [1] W. Bechtel and J. Mundale. Multiple realizability revisited: Linking cognitive and
349 neural states. *Philosophy of Science*, 66:175–207, 1999.
- 350 [2] N. Block and J. Fodor. What psychological states are not. *Philosophical Review*, 81,
351 1972.
- 352 [3] A. Church. An unsolvable problem of elementary number theory. *American Journal*
353 *of Mathematics*, 58:345–363, 1936.

- 354 [4] P.M. Churchland. Eliminative materialism and the propositional attitudes. *The Journal of Philosophy*, 78:67–90, 1981.
- 355
- 356 [5] P.M. Churchland. Functionalism at Forty: a critical retrospective. *Journal of Philosophy*, 1:33–50, 2005.
- 357
- 358 [6] J.B. Copeland. The Church-Turing thesis. *Stanford Encyclopaedia of Philosophy* [available online: <http://plato.stanford.edu/entries/church-turing/>], (last modified: August 2002, last consulted: July 2007).
- 359
- 360
- 361 [7] J. Fodor. Special sciences (or: on the disunity of science as a working hypothesis). *Synthese*, 28:97–115, 1974.
- 362
- 363 [8] Jerry Fodor. Special sciences: still autonomous after all these years. *Noûs*, 31:149–163, 1997.
- 364
- 365 [9] J. Kim. The myth of nonreductive materialism. *Proceedings and Addresses of the American Philosophical Association*, 63:31–47, 1989.
- 366
- 367 [10] J. Kim. Multiple realization and the metaphysics of reduction. *Philosophy and Phenomenological Research*, 52:1–26, (1992).
- 368
- 369 [11] H. Putnam, Psychological predicates. In *Art, Mind and Religion*, pages 37–48. University of Pittsburgh Press, Pittsburgh 1967.
- 370
- 371 [12] H. Putnam. *Mind, Language and Reality*. *Philosophical Papers* vol. 2. Cambridge University Press, Cambridge, 1975.
- 372
- 373 [13] A. Rosenberg. On multiple realization and the special sciences. *The Journal of Philosophy*, 98:365–373, 2001.
- 374
- 375 [14] S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*, Prentice-Hall, 1995.
- 376
- 377 [15] L. Shapiro. Multiple realization. *The Journal of Philosophy*, 97:635–654, 2000.
- 378 [16] A.M. Turing. On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, series 2, 42:230–265, 1936–37.
- 379
- 380 [17] N. Zangwill. Variable realization: Not proved. *The Philosophical Quarterly*, 42:214–219, 1992.
- 381