

A Characterization of Processing Loops in AI and Biological Systems and its Implications for Understanding Consciousness

1 **Malcolm J. Lett¹**

2 ¹No affiliation

3 *** Correspondence:**

4 Malcolm Lett

5 malcolm.lett@gmail.com

6 **Keywords: Access Consciousness, Phenomenal Consciousness, Biosemiotic Process, Visceral**
7 **Loop, Monitoring and Control, Recurrent Process, Non-physical Action, Conscious Content**

8 **Abstract**

9 The claim is made that recurrent cycles of non-physical "mental" actions are required in agents that
10 operate within complex environments, and that in order to maintain stability such actions require
11 regulation through the use of a model. The *Visceral Loop* is proposed as a novel framework for
12 characterizing three distinct kinds of processing within such a system in terms of how the system
13 uses that model. It is shown how this can be used to characterize human thought, including about
14 thought itself, and to explain the semiotic process when an individual concludes that they are
15 conscious. A proof is given relating an upper bound on data available within access consciousness to
16 the Visceral Loop characterizations of thought.

17 **1 Introduction**

18 Any computational system is limited in the complexity that it can handle within a single
19 computational step. For embodied agents, this appears as a limit on the environmental complexity
20 that they can sufficiently model and respond to within a single observe-infer-act cycle. For more
21 complex problems, multiple iterations of processing are required in order to determine the next
22 physical action. Such recurrency in processing may entail, for example, further analysis of the
23 environment in order to better model its state; or consideration of alternative action plans.

24 In biology, this provides scope for evolutionary pressures to trade off between a more energy hungry
25 complex brain and a simpler less energy intensive one that takes longer to make some decisions. Van
26 Bergen and Kriegeskorte (2020) make the case that recurrency is indeed employed in biology for that
27 very reason. A growing body of research in artificial intelligence is also now employing recurrency
28 and is showing that complex results can be achieved with shallower networks when using recurrency,
29 for example that of Kubilius et al. (2019) and Wen et al. (2018).

30 This paper makes the claim that an agent that employs *multi-step processing* (ie: multiple cycles of
31 processing without producing physical action) also must employ a model of its own processing
32 capabilities in order to regulate its *non-physical actions*. Different agentive architectures support
33 different abilities for the agent to introspect the internal structure of that model. In this paper the
34 *Visceral Loop*, is offered as novel framework for characterizing an agent's ability to introspect those
35 models and then use them for drawing inferences about itself.

The descriptive power of the Visceral Loop will be illustrated through three human-centric examples: i) how an individual may reach the conclusion that they are conscious, ii) providing an upper bound on the data content of consciousness, and iii) providing a new interpretation of neurobiological studies suggesting that reported awareness occurs after the act of decision making.

In the rest of this paper, the Regulation and Model sections elaborate further the need for regulatory models within any biological or artificial agent. The Schemas section suggests how this is manifested within humans in the form of a hypothesized *mind schema*. The Visceral Loop section presents the core thesis of this paper in the form of a mathematical formalism over the inferences that may be drawn from different kinds of model. The Consciousness section examines how the Visceral Loop can be applied to understanding aspects of consciousness through the presentation of the aforementioned humans. Final thoughts are presented in the Summary section.

2 Regulation

An autonomous embodied agent, depending on its purpose, may need to control either its environment, itself, or both, towards some static or dynamically determined target. That agent can be described as a *regulator* of its target system.

For example, an agent that regulates its environment operates within a system containing environment state s_{env} which changes with some ambient dynamics $d_{env}(t)$. The agent must perform an action, a_{env} , against the environment in order to regulate itself towards some target. After an action has been executed the environment state outcome o_{env} is influenced by both $d_{env}(t)$ and a_{env} . This can be summarized as the following equation:

$$s_{env} + d_{env}(t) + a_{env} = o_{env}$$

According to the *good regulator theorem*, if the agent is to regulate the environment state it must be a "model of the system" (Conant & Ashby, 1970). Furthermore, we can say that the efficiency of the agent to regulate its environment depends on its accuracy in modeling the system. Errors in the accuracy of the model result in errors in the regulation of the system. In learning agents, those errors are used for subsequent training of the model.

An embodied agent with complex actions may require an additional level of regulation. For example, an animal must not only regulate its external environment, but also regulate its own physical state. This includes both maintaining homeostasis and controlling the efficiency and effectiveness of its actions against the environment. The agent performs action a_{body} against its body with the intent to regulate the body towards efficiently achieving environment action a_{env} while satisfying its requirement for body homeostasis. Such an agent thus operates in a system that additionally has body state s_{body} with ambient dynamics $d_{body}(t)$. The agent performs action a_{body} against its body, producing outcome o_{body} , summarized as follows:

$$s_{body} + d_{body}(t) + a_{body} = o_{body}$$

Agents that incorporate multi-step processing have a third kind of action: one that changes its internal data state without affecting its physical state. This system requires regulation for the same reasons as for environment and body, but such *non-physical* actions do not elicit any change to s_{body} or s_{env} . Thus the agent must regulate its non-physical state s_{mind} , having ambient dynamics $d_{mind}(t)$. The target state in this case is dynamically inferred based on its requirement for environment action a_{env} , body action

76 a_{body} , and possibly for some form of non-physical homeostasis of s_{mind} . In order to regulate towards
 77 that target it performs action a_{mind} producing outcome o_{mind} , summarized as follows:

$$78 \quad s_{mind} + d_{mind}(t) + a_{mind} = o_{mind}$$

79 By way of example of the importance of such mind regulation, consider the case of fluent aphasia
 80 caused by damage to the Wernicke's area¹ of the brain. Individuals with fluent aphasia can easily
 81 produce speech, but it is typically full of many meaningless words and often unnecessarily long
 82 winded. Wernicke's area is associated with language comprehension. In a neurotypical individual, the
 83 comprehension of their own vocalizations provides a corrective mechanism. This illustrates the
 84 importance of feedback in the regulation of one's own actions, and by way of analogy extends to the
 85 regulation of non-physical actions.

86 3 Models

87 All of the systems described above are of the form $s + d(t) + a = o$. The production of the optimal
 88 action a for a given situation can be computed by a function, f , such that $a = f(s, o, t) = o - s - d(t)$. In
 89 this way, function f becomes a *model* of the system in exactly the way meant by Conant and Ashbey.
 90 There are many different ways of constructing such a function, with implications on how much its
 91 inherent model can be introspected for purposes other than merely computing the next action.

92 Consider the following function. This function is, for example, effective at predicting the action
 93 required to regulate towards a target state of 3 by doubling the input signal and comparing to that
 94 target state. However, an agent that merely uses this function to calculate actions cannot inspect
 95 anything about the function other than the actions it calculates for different inputs.

$$96 \quad f(x) = 3 - 2x$$

97 Alternatively consider Figure 1, which shows an abstract syntax tree² (AST) of the function above, of
 98 the sort used by computer science to parse an expression within a software compiler. Instead of using
 99 the above function, a regulating agent could use this AST to calculate its next action and achieve the
 100 same outcome. However the AST is a more explicit model of the dynamics being regulated. The
 101 components of the original function are represented individually and thus they can be individually
 102 queried. So here the AST can be introspected and much more can be derived from it that may apply
 103 either to the system being modeled or to how the AST models that system.

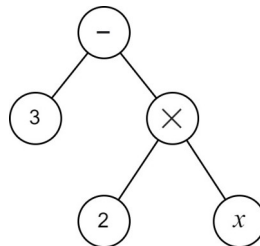


Figure 1: Abstract
syntax tree of $3 - 2x$

1 ¹https://en.wikipedia.org/wiki/Wernicke%27s_area

2 ²https://en.wikipedia.org/wiki/Abstract_syntax_tree

105 To examine the introspective opportunities further, consider the task of constructing a set, F , that
106 contains all beliefs that may be drawn from the model. In the case of the function, pairs of input and
107 output action values are all that can be drawn from the model, ie: $\langle 0, 3 \rangle$, $\langle 2, -1 \rangle$, $\langle -3, 9 \rangle$, $\langle -1, 5 \rangle$, etc.
108 The AST supports the ability to draw those same pairs of input and action values. However the AST
109 also supports that many other beliefs may be drawn from the model and added to F . For example that
110 i) the target is 3, ii) input signal x is significant to the calculation, while y and z are not, and iii) the
111 execution of the function depends on the operations of *subtraction* and *multiplication*.

112 So, it is clear that different architectures enable different levels of *introspection* of the underlying
113 models. What about the case for neural networks? In the modern use of artificial neural networks
114 (ANNs), it is commonplace to refer to ANNs as a *function approximator* (Goodfellow et al., 2016),
115 and indeed many networks fall into the category of a function. For example, in *model-free* deep
116 reinforcement learning (RL) an ANN is used to calculate either the next action or the expected value
117 of all possible actions given the current state (Lazaridis, Fachantidis & Vlahavas, 2020). The
118 architecture of the RL algorithm treats the ANN as a function without any introspective capabilities.
119 See Figure 2(A) for an example. There is also *model-based* RL. One variant of model-based RL,
120 illustrated in Figure 2(C), uses ANNs to predict the expected outcome of executing an action. The
121 introspective ability here is the same as for model-free deep RL - the ANN is treated as a function.
122 For the RL models mentioned so far, the set F of beliefs is of similar content: F is the set of
123 $\langle \text{state}, \text{action} \rangle$ or $\langle \text{state}, \text{action}, \text{outcome} \rangle$ tuples. There do exist forms of model-based RL that use
124 something more akin to the AST, usually where there is a known physics model that is represented
125 mathematically, and which could potentially be used to introspect for more than just
126 $\langle \text{state}, \text{action}, \text{outcome} \rangle$ tuples, such as is illustrated in Figure 2(B). However, a significant point to
127 note here is that ANNs, and probably neural networks in general, do not lend themselves to
128 introspection on their own.

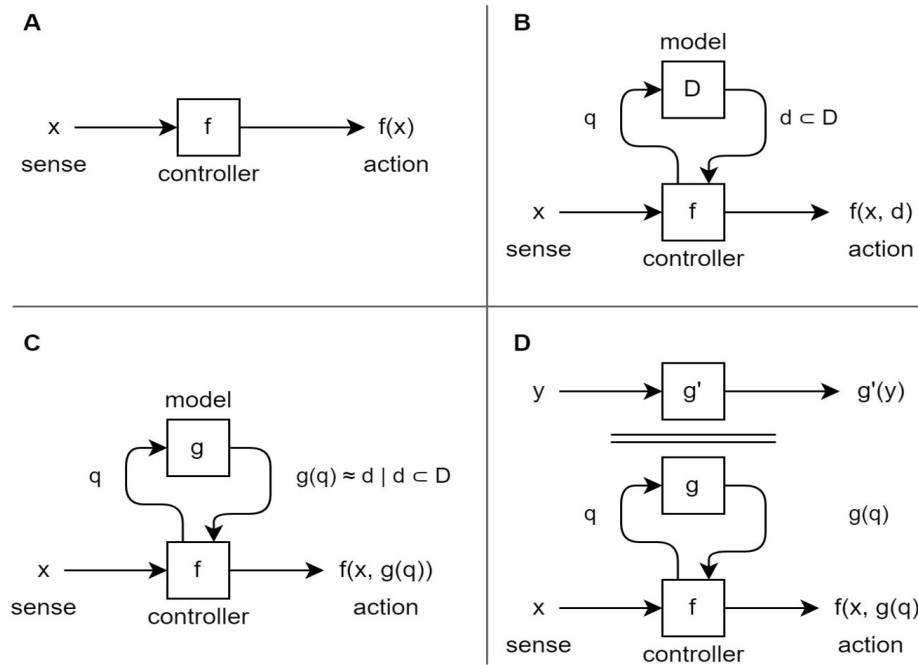


Figure 2: Different architectures for modeling regulatory actions against the environment. In (A), the controller determines the next action by executing function f against sense input x . The function may, for example, be an ANN that is trained through many iterations of an appropriate learning algorithm. Function f merely models the best next action without modeling any other aspects of that environment and thus cannot be used to introspect anything other than the next regulatory action. In (B), set D holds an explicit model of the environment which can be arbitrarily queried (q) to gain insight about any aspect of the environment that is encapsulated within D . Controller function f uses that to determine the next action. In (C), set D is replaced by function $g(q)$ which approximates queries against D . This architecture is commonly used in AI where the dynamics of the environment are too complex to determine a priori, and $g(q)$ is built as a second ANN that is trained through exploration. In (D), the secondary system $y = g'(y)$ models some aspect of the environment other than the next regulatory action. For example, it may observe and predict long term trends in the environment state. Potentially further additional modeling systems are required for each additional aspect of the environment that needs to be modeled.

130 For that reason, a third form of modeling system exists, whereby a secondary model predicts the
131 behaviors of the former, such as is illustrated in Figure 2(D). The secondary model may, for example,
132 be a second ANN that captures aspects of the same underlying system but at a more macro level, and
133 it may be more suitable for integration with other data. This macro representation is at the core of the
134 theory of Higher Order Thought Theory (HOTT) (Rosenthal, 1997 & 2006), and of recent theories
135 based on it such as Hierarchical Active Inference (Giovanni et al, 2018) and Integrated World
136 Modeling Theory (IWMT) (Safron, 2020).

137 4 Schemas

138 The lack of introspective ability of a simple function contrasts with the introspective ability of a
139 human. Psychology has long identified in humans the existence of a model of the individual's body –
140 known as the *body schema*. A good definition is given by Morasso et al (2015):

141 *"In summary, we view the body schema as a set of fronto-parietal networks that integrate*
142 *information originating from regions of the body and external space in a way, which is*
143 *functionally relevant to specific actions performed by different body parts. As such, the*
144 *body schema is a representation of the body's spatial properties, including the length of*
145 *limbs and limb segments, their arrangement, the configuration of the segments in space,*
146 *and the shape of the body surface".*

147 So the body schema is a model used in production of action control by integrating information from
148 our main physical senses and the proprioceptive senses (Proske & Gandevia, 2012). That model can
149 also be introspected – for example, we can know where our hands and feet are without seeing them –
150 and those introspections can become the topic of subsequent thought. But there are aspects of the
151 model that cannot be introspected – for example, we have no observability of the arrangement of the
152 sense nerves used to infer the hand and feet positions, or of the effector nerves used to actuate their
153 muscles.

154 This paper hypothesizes the existence of a second schema, the *mind schema*, that performs an
155 analogous role for the regulation of the mind and non-physical actions. Anecdotally, this seems
156 highly plausible within humans given our introspective ability towards our own mind's capabilities.
157 For example, we can know that we are good at focusing, but struggle with math, that we are more
158 creative when background music is present, and that we need the support of tools to help remember
159 people's names (eg: a notebook). The underlying notion here is that the mind schema helps us to
160 control, monitor, predict, and rationalize about our mental structure and actions in the same way that
161 our body schema does that for our physical structure and actions. It is the regulatory model for our
162 non-physical actions. Additionally, just as for the body schema, there is a distinct delineation
163 between what can be introspected and subsequently thought about, and what cannot.

164 The suggestion of a mind schema has also been made in the form of *Attention Schema Theory*
165 (Graziano & Kastner, 2011; Webb & Graziano, 2015; Graziano, 2017), although the meaning there is
166 perhaps narrower than what is proposed in this paper.

167 5 Visceral Loop

168 This paper introduces the *Visceral Loop* as a novel framework for the characterization of inferences
169 drawn by a processing loop within a biological or AI agent. The Visceral Loop is so named because it
170 refers to an agent concluding that it experiences consciousness in a visceral way. It identifies that an
171 agent with sufficient representational capabilities can, at the most optimum, conclude itself as
172 conscious within three iterations of the processing loop. Each of those iterations have specific
173 characteristics, and the Visceral Loop can be used to characterize any thought as falling into one of
174 those three iterations.

175 Let:

- 176 • E be the agent's set of beliefs about the external world

- 177 • B be the agent's set of beliefs about its own physical body (drawn from the body schema) and
178 of bodies in general, and if it has a concept of its identity then this set includes a belief that
179 relates other body beliefs to its identity
 - 180 • M be the agent's set of beliefs about its own mind (drawn from the mind schema) and of
181 minds in general, and if it has a concept of its identity then this set includes a belief that
182 relates other mind beliefs to its identity
 - 183 • $f(..)$ be the function executed by the agent on the specified inputs in order to draw inferences
- 184 M can be thought of as an agent's "theory of mind", because it relates not only to itself but also to its
185 ability to predict the hidden mental state of others.

186 5.1 Iteration 1

187 *Iteration 1* represents the most common kind of data processing, such as spending multiple
188 processing cycles to refine the identification of something within the visual field. While an agent's
189 mind schema may be used to regulate the thought process, the result of Iteration 1 never makes any
190 reference to it.

191 Let x be an inference produced as the result of a processing step, such that it does not draw any
192 reference to M (ie: $x \notin M$, and if x is a relation then $x = relation(\alpha, \beta)$ such that $\alpha \notin M$ and $\beta \notin M$ and
193 $\alpha \not\subset M$ and $\beta \not\subset M$). Given some sense input or past state s , a processing step is characterized as
194 Visceral Loop Iteration 1 if it is of the following form:

$$195 f(s, E \cup B \cup M) \rightarrow x$$

196 5.2 Iteration 2

197 *Iteration 2* processing steps draw conclusions that relate past non-physical actions and conclusions to
198 the agent's theory of mind and to the agent's concept of its identity. For example, concluding that a
199 past data state or non-physical action is classified as "thought", concluding whether the primary
200 source of a past data state was external or internal, or relating the fact of an internal source to the
201 agent's concept of its identity.

202 Iteration 2 requires an agent to have sufficient representational capabilities to produce inferences that
203 represent relations involving M . Given some prior inference y , a processing step is characterized as
204 Visceral Loop Iteration 2 if it is of the following form, and the relation with respect to M is non-
205 empty, and it can not be characterized as Iteration 3:

$$206 f(y, E \cup B \cup M) \rightarrow relation(y, M)$$

207 5.3 Iteration 3

208 *Iteration 3* is a special case of what would otherwise be Iteration 2, but it implies stricter
209 requirements on the agent's introspective and representational capabilities. Iteration 3 covers the
210 ability for the agent to develop a summary of its own mental capabilities (ie: some subset $m \subset M$),
211 and to consider that in relation to its conception of mental capabilities in general or to its identity (ie:
212 M). Iteration 3 is involved in an agent concluding itself as conscious, as will be seen in the section
213 below.

214 Given some prior inference $relation(z, M)$, and some subset of beliefs $m \subset M$, a processing step is
215 characterized as Visceral Loop Iteration 3 if it is of the following form and the relation with respect
216 to M is non-empty:

$$217 \quad f(relation(z, M), E \cup B \cup M) \rightarrow relation(m, M)$$

218 6 Consciousness

219 The Visceral Loop has implications for understanding consciousness; three examples of which will
220 be detailed shortly. For example, based on the Visceral Loop one can make an argument that
221 associates the behaviors of the underlying modeling mechanism to limits on the ability for self-
222 reference and to the ability for an individual to recognize their own consciousness. However, it is
223 necessary first to define what is meant here by the term *consciousness*.

224 Firstly, as so little is known of the nature of consciousness within non-human animals or in artificial
225 agents, the discussion will be limited to that of humans. Within that constraint, this paper focuses on
226 an individual's ability to have subjective experience of its observations (eg: of the external
227 environment), to be able to rationalize (a.k.a. "think") about things, and to have subjective experience
228 of those rationalizations. That domain of consciousness is typically split into two closely incident but
229 distinct forms:

- 230 • *Access consciousness* (A-Cs) refers to state representations that are "(1) inferentially
231 promiscuous, that is, poised for use as a premise in reasoning, (2) poised for rational control
232 of action, and (3) poised for rational control of speech." (Block, 1995). In other words, A-Cs
233 is the data content that has a direct influence on subsequent thought. For example, generally
234 the data content of access consciousness can be directly consciously thought about and
235 reported on via speech or other actions.
- 236 • *Phenomenal consciousness* (P-Cs), in contrast, refers to subjective experience and to the
237 question of why a physical bundle of matter should have such a thing. It is often described as
238 "what it is like to be" (Nagel, 1974) a conscious organism. In approximate terms, P-Cs is the
239 phenomenal aspect of having a subjective experience of the data content of A-Cs.

240 The discussions in this paper refer to the data content of both A-Cs and P-Cs, but in general does not
241 draw an explicit distinction between A-Cs versus P-Cs, under the assumption that there is a single
242 shared physical mechanism underlying them both. Here, *data content* is used merely to draw a
243 distinction between what is and what is not present within A-Cs or P-Cs. It is not intended to imply
244 any particular representational structure. For example, whether or not someone can see and report on
245 something within their visual field determines part of the data content of A-Cs, and whether or not
246 someone phenomenally experiences pain determines part of the data content of P-Cs.

247 6.1 Concluding oneself as conscious

248 In this first example, the Visceral Loop is applied to understand the thought processes whereby an
249 individual concludes themselves as conscious.

250 Consider the following sequence of internal mental observations:

- 251 1. "What's that red blob in the tree? Oh, it's an apple".

2. "Those thoughts just came from my mind, and not from the outside world".

3. "That's what consciousness is. I am conscious".

The first observation is a straightforward example of Visceral Loop Iteration 1 that does not make any reference to the agent's theory of mind (of their own mind or of others). With reference to the formal definition of the Visceral Loop in the prior chapter, the concepts of "red", "blob", "tree" and "apple" are all contained within the set E , and thus the inference in relation to the visual field sense input s is of the form $x_1 = \text{relation}(s, E)$.

The second observation contains two examples of Iteration 2 inferences. In the first, the individual's processing capabilities have selected attentional focus upon the prior Iteration 1 inference, and have drawn a subsequent inference about it as being data that can be classified as a "thought". As beliefs about "thought" are contained within M , this is an inference of the form $x_2 = \text{relation}(x_1, M)$. In the second, the individual draws a subsequent inference about the source of the Iteration 1 inference as being their own mind. The individual's ability to classify inferences in relation to themselves also depends upon M , and the inference is of the form $x_3 = \text{relation}(x_1, M)$.

The third observation draws upon the individual having an a priori conception about consciousness in general, denoted here by $m_c \subset M$. The individual compares its prior Iteration 2 inferences x_2 and x_3 to m_c , and produces an inference that x_2 and x_3 together satisfy the requirements for consciousness. This is another iteration 2 inference of the form $x_4 = \text{relation}(x_2 \wedge x_3, m_c)$. Finally, the individual relates m_c , the belief of consciousness in general, to itself, which again depends on M . That final inference is thus an Iteration 3 inference in the form $x_5 = \text{relation}(m_c, M)$.

6.2 Content of conscious thought

As a second example of the descriptive power of the Visceral Loop, a theorem is now presented that makes the claim that the Visceral Loop explains the data content of conscious experience. An exact distinction of how it applies to A-Cs versus P-Cs is omitted here, with some discussion of its issues included in a later section.

First an axiomatic baseline must be established. It seems reasonable to expect that there is no way in which an individual may consciously experience something and yet be unable to subsequently think about that experience and to know that they are thinking about that experience. Indeed, this is consistent with the *transitivity principle* of Rosenthal (1997). Thus, it would seem that being able to knowingly think about our conscious experiences is a fundamental component of consciousness. The following claims are derived from this statement without further proof:

Claim 1:

- All conscious experience is subsequently available for further thought.

Claim 2:

- For all thought about conscious experience, the individual can identify that thought as being their own.

288 Note that these claims do not assume that all conscious experience is subsequently thought about;
289 only that it is in principle available for such thought.

290 Theorem 1:

- 291 • the data content of conscious experience is upper bounded by the data about which Visceral
292 Loop iteration 2 inferences can be produced.

293 Proof:

- 294 • Thought is a computational process, and thus is a series of inferences.
- 295 • As per claim 1, all of conscious experience must be available for producing subsequent
296 inferences about those conscious experiences.
- 297 • As per claim 2, the individual must be able to identify that they produced those inferences.
- 298 • In order for an individual to identify an inference as being their own, they must have some
299 beliefs about their inference capabilities and how they relate to themselves as an individual
300 entity. This is included in the set M , which iteration 2 produces inferences in relation to, and
301 which is not directly accessible for inferences within iteration 1.
- 302 • Imagine some supposed experience, and an inference i produced about that experience.
303 Additionally imagine that an iteration 2 inference cannot be produced about i , for example,
304 due to some incompatibility of structure, lack of data path to iteration 2 processing
305 capabilities, or inherent limitation in iteration 2 processing capabilities. The inference i cannot
306 be identified in relation to the individual. As such, the supposed experience fails on Claim 2
307 and i must be in actual fact an inference about some sort of non-conscious experience.
- 308 • Thus, an experience is not a conscious experience if it can only lead to inferences which
309 cannot be included in an iteration 2 inference.

310 6.3 Delayed awareness of decisions

311 The Visceral Loop can be used to understand other aspects of thought. For example, fMRI and EEG
312 studies have suggested that we become aware of a decision after it is made (Soon et al, 2008; Libet et
313 al, 1983). At first glance this might seem to suggest that our conscious thought is non-causative,
314 instead being just some sort of after-the-fact passive summarization, such as is claimed by the theory
315 of Epiphenomenalism³.

316 The framework of the Visceral Loop provides a different interpretation. It explains that the act of
317 making a decision and the conscious consideration of having made that decision occur in different
318 processing cycles. First one or more processing cycles are used to reach the decision and to trigger
319 the resultant action. Subsequently the individual may use one or more additional processing cycles to
320 examine their most recent inference. So it is entirely predictable that the individual would activate
321 brain regions for reporting the decision after a measurable delay from the decision itself being made
322 and a resultant action initiated. Importantly, the same underlying system can produce both sets of

3 ³<https://plato.stanford.edu/entries/epiphenomenalism>

323 processing cycles, and thus there is no reason to conclude from fMRI or EEG delays that
324 consciousness is non-causal.

325 In short, we can only think about one thing at a time, so the decision itself and thought about the
326 decision require separate processing steps.

327 7 Analysis

328 The examples above show that it is possible to mathematically reason about thought processes, their
329 sequencing, and what can and cannot be thought about depending on the capabilities of the
330 underlying thought processing system. This section examines the Visceral Loop against some
331 existing theories and points out some existing weaknesses.

332 7.1 Visceral Loop as a Biosemiotic Process

333 The Visceral Loop views thought as a *biosemiotic process*⁴. Semiotics is the study of signs and their
334 interpretations, and biosemiotics looks at how semiotics plays out within biology, including neural
335 cognition. A semiotic process has three components: a *referent*, the object for which a sign or
336 representation will be made; a *representamen*, which is a representation of the referent, a.k.a. a
337 "sign"; and a *interpretant*, the interpretation made from the representamen. The interpretant may or
338 may not accurately reflect the original referent, depending on the quality of representamen and on the
339 ability of the system (ie: the interpreter) to infer information about the referent from the
340 representamen. For example, in Figure 1 the referent of someone liking something is converted into a
341 thumbs-up representamen, but it has multiple possible interpretants depending on the context and
342 background of the interpreter.



Figure 3: An example semiotic process. Someone wishes to indicate that they like something - their liking of that something is the *referent* that they wish to convey through a sign. The sign they use, the *representamen*, is the icon of a hand with thumb pointing up. The representamen is then interpreted by someone else, forming the *interpretant* within the mind of that second individual. Depending on the context and the cultural background of the second individual, they may form one of multiple possible interpretants (eg: a thumbs up gesture is considered an insult in some cultures).

344 From an external point of view, the brain acts as a semiotic process against the environment. As
345 illustrated in Figure 4, the environment is the referent, with a concrete and actual state that the brain
346 does not have direct awareness of. Instead, the brain uses the physical senses as a representamen of
347 that environment, and from that representation produces an interpretant.

4 <https://en.wikipedia.org/wiki/Biosemitics>

348 The Visceral Loop is interested in a more internal semiotic process within the brain. As illustrated in
349 Figure 5, the interpretant from before must itself be encoded in a representation. And that
350 representation is subject to the constraints imposed by the specific characteristics of the underlying
351 biological substrate. In order for the brain to subsequently use that representation, it must decode its
352 meaning. Thus an inner semiotic process is activated, where the prior inference becomes the referent,
353 which is encoded as a representamen, and subsequently decoded as a new interpretant, possibly after
354 combining with additional information from other sources.

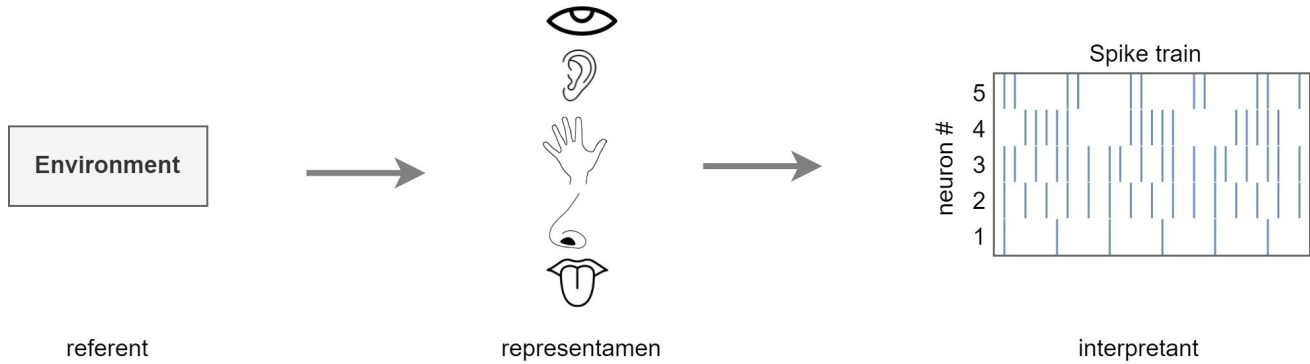


Figure 4: Semiotics of environment interpretation

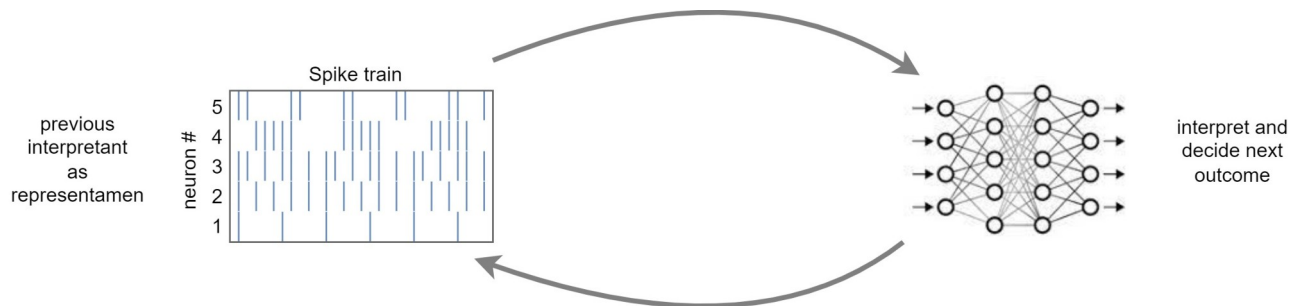


Figure 5: Semiotics of internal thought

357 The semiotic interpretation of cognition is important because it explains that a data state is
358 meaningless without a process that interprets that data state; even if the process generated the data
359 state in the first place, that data state is still meaningless to the process without re-interpreting it.
360 Thus, while the brain must include considerable machinery for the interpretation of external stimuli,
361 it must also include considerable machinery for the interpretation of its own outputs. The Visceral
362 Loop builds on this by describing the kinds of interpretations that may be generated depending on the
363 what information is available within the representamen and on the capabilities of the process.

364 7.2 Relationship to Higher Order Thought Theory and Global Workspace Theory

365 The Visceral Loop and the biosemiotic process that it operates upon examine thought at a high-level,
366 but the Visceral Loop does not presuppose a particular physical or computational structure. In that
367 sense it is applicable to many theories of cognition and brain function. However, it is easiest to see its
368 consistency with other theories that also examine thought at a high-level or that provide theories

369 about the underlying mechanisms of such high-level thought. For example, those of Higher-Order
370 Thought Theory (HOTT), and Global Workspace Theory (GWT).

371 HOTT describes mental processes as hierarchical (Rosenthal 1997; Rosenthal 2006). In HOTT,
372 consciousness forms the highest layer of the hierarchy, and has access to the output from the
373 immediate layer below, but not of further lower layers. Thus the information available for conscious
374 thought is a unified high level abstract representation formed from the output of many sub-processes.
375 The Visceral Loop does not require a HOTT architecture, but it explains how a HOTT architecture
376 results in the specific nature of human experience in terms of their being a distinction between those
377 internal processes which are directly observable, versus those which are not.

378 The Visceral Loop can also be seen as a monitoring and control device that governs the actions of the
379 rest of the agent's cognitive function. This suggests why a higher-order layer is requisite. The internal
380 processes of the entire brain are too complex to monitor at their native low-level neuronal states. The
381 global monitoring and control function of the brain needs to know only the broad trends in behavior
382 in order to predict whether the existing behavior is suitable for the larger objective. Thus it is more
383 efficient that it operates over a higher-order unified representation that contains only the information
384 that is pertinent for the purpose of self-governing.

385 Conveniently, this finally explains why we are conscious of certain processes within our own brains
386 and not conscious of others – that our brains create a higher-order representation for the purpose of
387 self-governing, and only that self-governing function has the full semiotic process capable of
388 producing conscious experience, and thus we are only conscious of the data content that is presented
389 to that self-governing process.

390 **7.3 A step towards understanding General Intelligence**

391 The problem of solving Artificial General Intelligence (AGI) has received renewed interest of late
392 (Adams et al., 2012; van Gerven, 2017; Cervantes et al., 2021). The Visceral Loop may provide some
393 insight.

394 Firstly, an AGI operates within a complex environment, has complex responses, and thus it needs to
395 incorporate recurrent non-physical actions, as per the arguments presented within this paper. This
396 implies the need for the agent to introspect and monitor its own mental behaviors in order to control
397 them for stability and efficiency. That monitoring and control operates as a semiotic process against
398 the agent's own internal mental representations. Thus an AGI should be expected to exhibit the
399 characteristics that are measured by the Visceral Loop, and one aspect of the level of advancement of
400 the AGI can be measured by whether it can execute Visceral Loop Iterations 2 and 3.

401 This suggests a particular AGI architecture whereby a discrete sub-component provides monitoring
402 and control and governs the actions of multiple other sub-components, each with individual
403 specializations, such as is illustrated in Figure 6. Were those sub-components to operate against each
404 other through competitive processes alone, they may be at risk of chaotic behavior and runaway
405 instabilities. Provided with an appropriate learning algorithm and an appropriate objective that
406 measures the agents overall stability and efficiency, an aggregate system that incorporates a global
407 monitoring and control function can efficiently learn to govern itself and to avoid chaotic runaway
408 instabilities. For example, one such learning algorithm and objective is that of Active Inference
409 (Friston, Daunizeau, and Kiebel, 2009; Tschantz, Baltieri, et al, 2020; Sajid, Ball, et al, 2021).

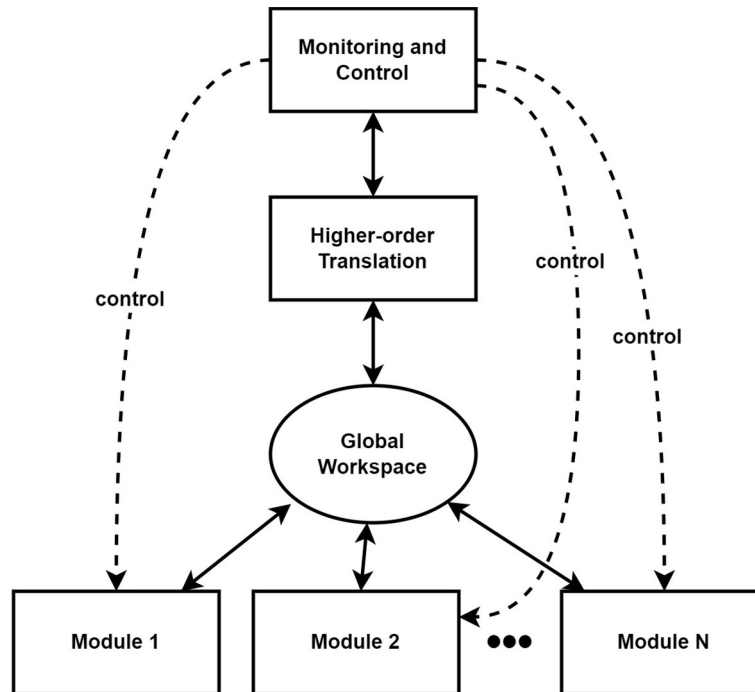


Figure 6: A potential AGI architecture. A Monitoring and Control unit governs the execution of individual specialized modules. Based on discussions within Section 7.2, the monitoring and control unit views the behaviors of the individual modules through a higher-order translation, and applies control either directly or through reverse higher-order translation. A global workspace model is naively assumed for sharing of information between modules.

7.4 A step towards understanding Consciousness

Within section 6 the discussion was restricted to human consciousness, but the Visceral Loop provides a mathematical basis for understanding the potential for consciousness in other agents, including artificial ones. Where an agent has a processing architecture sufficient for producing Visceral Loop Iteration 2 inferences, then it has the potential for conscious thought.

The recurrent nature of the semiotic process underlying the Visceral Loop explains the need for a computational architecture that incorporates i) recurrency at the high-level scale, ii) persistence of state between recurrent cycles, and iii) representational integration of diverse computational outcomes.

Another similar theory that attempts to mathematically reason about levels of consciousness is Integrated Information Theory (IIT) (Tononi, 2004; Oizumi et al., 2014). IIT provides a mechanism for quantifying the amount of recurrent integration in a system, symbolized in its amount of Φ . Importantly, IIT is applicable to both biological and artificial agents. The Visceral Loop is thus compatible with and enhances IIT, in that it gives a higher-level reason for the need for recurrency.

Phenomenal experience is generally considered to be the hallmark of consciousness, making the distinction between an agent that is truly conscious versus a philosophical zombie of the same. The theory presented here offers some insight into the mechanism behind the ability for an individual to examine their own thought. This is assumed to be closely related to the content of the individual's phenomenal experience. Thus, the Visceral Loop is an important aspect of consciousness, and understanding of it may help lead towards artificial consciousness.

7.5 Limitations

The Visceral Loop offers significant insight into consciousness, but with the introduction of any new theory there are limitations and areas which may need improvement. It is thus worth briefly discussing some of the limitations with the Visceral Loop as stated, and with the examples of its use given above.

Theory 1, presented above, is weakened somewhat by the fact that it avoids clarifying whether it applies to A-Cs, P-Cs, or both. Unfortunately such an analysis is hampered due to presently irreconcilable differences in opinion about the relative natures and definitions of A-Cs and P-Cs, and about whether A-Cs or P-Cs can exist without the other⁵. For example, some authors make the claim that P-Cs can exist without A-Cs (Block, 1995; Armstrong, 1995). This would imply some kind of representational content within P-Cs that is independent of A-Cs and thus cannot be subsequently thought about (ie: violating Claim 1). There is also research suggesting that some A-Cs information can influence our subsequent thought without us having P-Cs of it, such as in the example of blind-sight (Block, 1995).

In both cases there remains heated debate with no obvious end in sight. The best we can say is that A-Cs and P-Cs largely correlate, although there may be exceptions. The author believes that stronger claims can yet be made about the natures of both A-Cs and P-Cs in terms of the Visceral Loop, but that shall remain the topic of a separate investigation.

The theories presented in this paper make no claim over the so called "hard problem of consciousness" – the question of why we have phenomenal experience at all (Chalmers, 1995). That remains an open question. The definition of the Visceral Loop itself is compatible with many of the existing arguments that attempt to address the hard problem, although Theory 1 is challenged by some such arguments.

One area in which the Visceral Loop does not suffer from a common limitation is that it avoids any risk of a recursive definition of conscious thought. Based on the Visceral Loop, a tentative definition of consciousness might be that:

We are conscious of a perception or thought if that perception or thought is in principle capable of being subsequently thought about via Visceral Loop Iterations 2 or 3, and we have subjective experience of it at the moment that an Iteration 2 or 3 result is produced.

An important point here is that the perception or thought does not have to be immediately processed by Iteration 2 or 3; it can pass via memory and be processed by Iteration 2 or 3 at a later moment. In fact, anything that is accessible from memory and which is capable of being processed through

⁵For a good example of the issues with defining A-Cs and P-Cs, see the responses to Block (1995), accessible at https://www.nyu.edu/gsas/dept/philo/faculty/block/papers/1995_Function.pdf

Iteration 2 or 3 (subject to the availability of the necessary data pathways) can be subjectively experienced.

8 Summary

This paper makes the claim that agents operating within complex environments and with complex bodies require a trade-off between the inferential computing power of a single processing step versus the use of multiple iterations of a recurrent processing loop. Processing loop iterations that do not lead to physical actions thus have non-physical actions, and in order to be a good regulator of those non-physical actions, the agent must model its non-physical behaviors.

It has been shown that some architectures for non-physical behavioral modeling enable introspective abilities, such as that enjoyed by humans. The Visceral Loop provides a mechanism to mathematically reason about those introspective abilities, and to classify the capabilities of different architectures. Where Integrated Information Theory (IIT) provides a mechanism for quantifying the amount of recurrent integration in a system, the theory behind the Visceral Loop explains why that recurrency is important in the first place. The Visceral Loop provides a mechanism for qualifying the introspective capabilities in a way that is simultaneously at a coarser grain and simpler to calculate than what IIT targets.

The Visceral Loop provides insight into consciousness in general. It shows that the content of consciousness is defined by a biosemiotic process that operates against a higher-order summary of the internal processes of the rest of the brain, and that it is necessary for stable functioning of a complex cognitive system. It also suggests a pathway towards building an Artificial General Intelligence. It suggests the introduction of a self-monitoring and control process that governs the execution of multiple specialized modules which.

9 Author Contributions

The author confirms being the sole contributor of this work and has approved it for publication.

10 Funding

No funding was received.

11 References

- Adams, S., Arel, I., Bach, J., Coop, R., Furlan, R., Goertzel, B., et al. (2012). Mapping the Landscape of Human-Level Artificial General Intelligence. *AI Magazine*, 33:1, 25-42, doi:10.1609/aimag.v33i1.2322.
- Armstrong, D. M. (1995). Perception-consciousness and action-consciousness? *Behav. Brain Sci.*, 18:2, 247-248, doi:10.1017/S0140525X0003819X.
- Block, N. (1995). On a confusion about a function of consciousness. *Behav. Brain Sci.*, 18:2, 227–247. doi:10.1017/S0140525X00038188.
- Cervantes, J. P., Martin, L., Dounce, I. A., Avila-Contreras, C., and Corchado, F. F. (2021). Methodological aspects for cognitive architectures construction: a study and proposal. *Art. Intel. Rev.*, 54, 2133-2192. doi:10.1007/S10462-020-09901-X.

- 500 Chalmers, D. J. (1995). Facing up to the problem of consciousness. *J. Conscious. Stud.*, 2:3, 200-19.
501 doi:10.1093/acprof:oso/9780195311105.003.0001.
- 502 Colagrosso, M. D., and Mozer, M. C. (2004). Theories Of Access Consciousness. *NeurNIPS* 2004.
- 503 Conant, R. C., and Ashby, W. R. (1970). Every good regulator of a system must be a model of that
504 system. *Int. J. Systems Sci.*, 1:2, 89-97. doi:10.1080/00207727008920220.
- 505 Giovanni, P., Rigoli, F., and Friston, K. J. (2018). Hierarchical Active Inference: A Theory of
506 Motivated Control. *Trends in Cognitive Sciences*, 2:4, pp 294-306. doi:10.1016/j.tics.2018.01.009.
- 507 Friston, K. J., Daunizeau, J., Kiebel, S. J. (2009). Reinforcement Learning or Active Inference?
508 PLOS ONE, 4:7, doi:10.1371/journal.pone.0006421.
- 509 Goodfellow, I., Bengio, Y., and Courville, A. (2016). Deep learning. MIT Press.
- 510 Graziano, M. S. A., and Kastner, S. (2011). Human consciousness and its relationship to social
511 neuroscience: a novel hypothesis. *Cogn. Neurosci.* 2, 98–113. doi:10.1080/17588928.2011.565121
- 512 Graziano, M. S. A. (2017). The Attention Schema Theory: A Foundation for Engineering Artificial
513 Consciousnes. *Front. Robot. AI.* doi:10.3389/frobt.2017.00060.
- 514 Kubilius, J., Schrimpf, M., Kar, K., Rajalingham, R., Hong, H., Majaj, N., Issa, E., et al. (2019).
515 Brain-like object recognition with high-performing shallow recurrent anns. *NeurIPS* 2019, 12805-
516 12816.
- 517 Lazaridis, A., Fachantidis, A., and Vlahavas, I. (2020). Deep Reinforcement Learning: A State-of-
518 the-Art Walkthrough. *J. Artif. Intell. Res.*, 69, pp 1421-1471. doi:10.1613/jair.1.12412.
- 519 Libet, B., Gleason, C. A., Wright, E. W., and Pearl, D. K. (1983). Time of conscious intention to act
520 in relation to onset of cerebral activity (readiness-potential). The unconscious initiation of a freely
521 voluntary act. *Brain: a journal of neurology*, 106:3, 623–642. doi:10.1093/brain/106.3.623.
- 522 Morasso, P., Casadio, M., Mohan, V., Rea, F., and Zenzeri, J. (2015). Revisiting the body-schema
523 concept in the context of whole-body postural-focal dynamics. *Front. Hum. Neurosci.* 9:83.
524 doi:10.3389/fnhum.2015.00083.
- 525 Nagel, T. (1974). What Is It Like to Be a Bat?. *The Philosophical Review*, 83:4, 435–450.
526 doi:10.2307/2183914.
- 527 Oizumi, M., Albantakis, L., and Tononi, G. (2014). From the phenomenology to the mechanisms of
528 consciousness: integrated information theory 3.0. *PLoS Comput. Biol.* 10:e1003588.
529 doi:10.1371/journal.pcbi.1003588.
- 530 Proske, U., and Gandevia, S. C. (2012). The Proprioceptive Senses: Their Roles in Signaling Body
531 Shape, Body Position and Movement, and Muscle Force. *Physiol. Rev.*, 92:4, 1651-1697.
532 doi:10.1152/physrev.00048.2011.
- 533 Rosenthal, D. M. (1997). "A Theory of Consciousness," in *The Nature of Consciousness:*
534 *Philosophical Debates*, ed. N. Block, O. Flanagan, and G. Güzeldere (Cambridge, MA: MIT
535 Press/Bradford Books), 729-753.

- 536 Rosenthal, D. M. (2006). "Consciousness and Higher-Order Thought," in Encyclopedia of Cognitive
537 Science, ed. L. Nadel, doi:10.1002/0470018860.s00149.
- 538 Safron A. (2020). An Integrated World Modeling Theory (IWMT) of Consciousness: Combining
539 Integrated Information and Global Neuronal Workspace Theories With the Free Energy Principle and
540 Active Inference Framework; Toward Solving the Hard Problem and Characterizing Agentic
541 Causation. *Front. Artif. Intell.*, 3:30, doi:10.3389/frai.2020.00030.
- 542 Sajid, N., Ball, P. J., Parr, T., Friston, K. J. (2021). Active Inference: Demystified and Compared.
543 *Neural. Comput.*, 33:3, 674–712, doi:10.1162/neco_a_01357.
- 544 Seth, A. K., Suzuki, K., and Critchley, H. D. (2012). An interoceptive predictive coding model of
545 conscious presence. *Front. Psychol.*, 2:395. doi:10.3389/fpsyg.2011.00395.
- 546 Soon, C. S., Brass, M., Heinze, H. J., & Haynes, J. D. (2008). Unconscious determinants of free
547 decisions in the human brain. *Nat. Neurosci.*, 11:5, 543–545, doi:10.1038/nn.2112.
- 548 Tschantz, A., Baltieri, M., Seth, A., and Buckley, C. (2020). Scaling Active Inference. *IJCNN*, 1-8,
549 doi: 10.1109/IJCNN48605.2020.9207382.
- 550 Tononi, G. (2004). An Information Integration Theory of Consciousness. *BMC Neurosci.*, 5:42,
551 doi:10.1186/1471-2202-5-42.
- 552 van Bergen, R. S., and Kriegeskorte, N. (2020). Going in circles is the way forward: the role of
553 recurrence in visual inference, *Curr. Opin. Neurobiol.*, 65, 176-193, doi:10.1016/j.conb.2020.11.009.
- 554 van Gerven, M. A. (2017). Computational Foundations of Natural Intelligence. *Front. Comp.*
555 *Neurosci.*, 11. doi:10.3389/fncom.2017.00112.
- 556 Webb, T. W., and Graziano, M. S. A. (2015). The attention schema theory: a mechanistic account of
557 subjective awareness. *Front. Psychol.* 6:500. doi:10.3389/fpsyg.2015.00500.
- 558 Wen, H., Han, K., Shi, J., Zhang, Y., Culurciello, E. and Liu, Z. (2018). Deep Predictive Coding
559 Network for Object Recognition. *Proceedings of the 35th International Conference on Machine*
560 *Learning*, in *Proc. Mach. Learn. Res.*, 80, 5266-5275