

# AI Consciousness is Inevitable: A Theoretical Computer Science Perspective

Lenore Blum and Manuel Blum

## ABSTRACT

We look at consciousness through the lens of Theoretical Computer Science, a branch of mathematics that studies computation under resource limitations, distinguishing functions that are efficiently computable from those that are not. From this perspective, we develop a formal machine model for consciousness. The model is inspired by Alan Turing's simple yet powerful model of computation and Bernard Baars' theater model of consciousness. Though extremely simple, the model aligns at a high level with many of the major scientific theories of human and animal consciousness, supporting our claim that machine consciousness is inevitable.

## Table of Contents

<b>1</b>	<b>INTRODUCTION .....</b>	<b>2</b>
<b>2</b>	<b>BRIEF OVERVIEW OF RCTM, A ROBOT WITH A CTM BRAIN.....</b>	<b>5</b>
2.1	FORMAL DEFINITIONS OF CTM/RCTM .....	5
2.2	CONSCIOUS ATTENTION IN RCTM .....	9
2.3	CONSCIOUS AWARENESS AND THE FEELING OF CONSCIOUSNESS IN RCTM .....	10
2.3.1	<i>The Model-of-the-World and Brainish co-Evolve .....</i>	<i>12</i>
2.3.2	<i>Conscious Awareness in rCTM .....</i>	<i>15</i>
2.4	RCTM AS A FRAMEWORK FOR ARTIFICIAL GENERAL INTELLIGENCE (AGI) .....	16
<b>3</b>	<b>ALIGNMENT OF RCTM WITH OTHER THEORIES OF CONSCIOUSNESS .....</b>	<b>17</b>
3.1	GLOBAL WORKSPACE (GW)/GLOBAL NEURONAL WORKSPACE (GNW) .....	17
3.2	ATTENTION SCHEMA THEORY (AST) .....	17
3.3	PREDICTIVE PROCESSING (PP) .....	17
3.4	EMBODIED EMBEDDED ENACTIVE EXTENDED MIND (EEEE MIND).....	18
3.5	INTEGRATED INFORMATION THEORY (IIT) .....	19
3.6	EVOLUTIONARY THEORIES OF CONSCIOUSNESS .....	19
3.7	EXTENDED RETICULOTHALAMIC ACTIVATING SYSTEM (ERTAS) + FREE ENERGY PRINCIPLE (FEP) .....	20
<b>4</b>	<b>ADDRESSING KEVIN MITCHELL'S QUESTIONS FROM THE PERSPECTIVE OF RCTM.....</b>	<b>21</b>
<b>5</b>	<b>SUMMARY AND CONCLUSIONS.....</b>	<b>32</b>
<b>6</b>	<b>ACKNOWLEDGEMENTS .....</b>	<b>33</b>
<b>7</b>	<b>APPENDIX.....</b>	<b>34</b>
7.1	A BRIEF HISTORY AND DESCRIPTION OF THE TCS APPROACH TO COMPUTATION .....	34
7.2	THE PROBABILISTIC COMPETITION FOR CONSCIOUS ATTENTION AND THE INFLUENCE OF DISPOSITION ON IT .....	35
	<b>REFERENCES .....</b>	<b>38</b>

## 1 Introduction

We study consciousness from the perspective of Theoretical Computer Science (TCS), a branch of mathematics concerned with understanding the underlying principles of computation and complexity, in particular, the implications and surprising consequences of resource limitations.

By taking resource limitations into account, the TCS perspective is distinguished from Turing's Theory of Computation (TOC) where limitations of time and space do not figure. TOC distinguishes computable from not computable. It does not distinguish between efficiently computable and not efficiently computable.<sup>1</sup> In our study of consciousness from a TCS perspective we highlight the importance of resource limitations for tackling consciousness and related topics such as the paradox of free will (Blum & Blum, 2022).

Elsewhere (Blum & Blum, 2021), we describe the Conscious Turing Machine (CTM), a *simple formal machine* model of *consciousness* inspired in part by Alan Turing's *simple, yet powerful, formal machine* model of *computation* (Turing, 1937), and by Bernard Baars' theater model of consciousness (Baars, Bernard J., 1997) and (Baars, 1997). In (Blum & Blum, 2022), we consider how a CTM could exhibit various phenomena associated with consciousness (e.g., blindsight, inattentive blindness, change blindness) and present CTM explanations that agree, at a high level, with cognitive neuroscience literature. In the spirit of TCS, informal notions are defined formally in the CTM model.

In contrast to Turing, we take resource limitations into account, both in designing the CTM model (e.g., size of chunks, speed of computation) and in how resource limitations affect (and help explain) phenomena related to consciousness (e.g., change blindness and dreams) and feelings of consciousness. Our perspective differs even more. What gives the CTM its feeling of consciousness is not its input-output map, nor its computing power, but what's under the hood.<sup>2</sup>

In this paper we take a brief look under the hood (section 2.1). We also indicate how the CTM provides a framework for constructing an Artificial General Intelligence (AGI). (See section 0.)

---

<sup>1</sup> For a brief history of TOC and TCS, see Appendix 7.1.

<sup>2</sup> This is important. We claim that simulations that modify CTM's key internal structures and processes will not in general experience what CTM experiences. On the other hand, we are not claiming that the CTM is the only possible machine model to experience feelings of consciousness. The CTM is a minimal machine model for consciousness.

## AI Consciousness is Inevitable

In addition, we show how the CTM naturally *aligns* with and *integrates* features considered key to human and animal consciousness by many of the major scientific theories of consciousness (Section 3).<sup>3</sup> These theories consider different aspects of consciousness and often compete with each other (Lenharo, 2024). Yet their alignment with the CTM at a high level helps demonstrate their compatibility and/or complementarity.

*Even more, their alignment with the CTM, a simple machine model that exhibits many of the important phenomena associated with consciousness, supports our claim that a conscious AI is inevitable.*<sup>4</sup>

David Chalmers' introduction of the Hard Problem (Chalmers, 1995) helped classify most notions of consciousness into one of two types. The first type, variously called access consciousness (Block, 1995) or computational consciousness or cognitive consciousness (Humphrey, 2023), we call *conscious attention* (section 2.2). The second type (associated with the Hard Problem) is called *subjective* or *phenomenal* consciousness and is generally associated with feelings or qualia. We call it *conscious awareness* (section 2.3). Chalmers' Hard Problem can be viewed as a challenge to show that subjective consciousness is "computational".

We argue (in our previously cited papers and in section 2.3) that consciousness generally requires both conscious attention and conscious awareness, each informing the other to various degrees. We contend that a machine that interacts with its worlds (inner and outer) via input sensors and output actuators, that constructs models of these worlds enabling planning, prediction, testing, and learning from feedback, and that develops a rich internal multimodal language, can have both types of consciousness. In particular, we contend that subjective consciousness is computational.

We emphasize that the CTM is a simple formal machine model designed to explore and understand consciousness from a TCS perspective. It is not intended to model the brain nor the neural

---

<sup>3</sup>These theories include: The Global Workspace/Global Neuronal Workspace (GW/GNW), Attention Schema Theory (AST), Predictive Processing (PP), Integrated Information Theory (IIT), Embodied, Embedded, Enacted and Extended (EEEE) theories, Evolutionary theories, and the ERTAS (Extended Reticulothalamic Activating System) + FEP (Free Energy Principle) theories.

Indeed, Wanja Wiese describes the CTM as a "minimal unifying model" (Wiese, 2020) and (personal communication).

<sup>4</sup> Thus, our response to the query, "could machines have it [consciousness]?" (Dehaene, Lau, & Kouider, 2017), is "yes". In a recent preprint, (Farisco, Evers, & Changeux, 2024) consider the question "Is artificial consciousness achievable?" from the perspective of the human brain. They conclude with a tentative possibility for "non human-like forms of consciousness."

## AI Consciousness is Inevitable

correlates of consciousness, though it is inspired by cognitive and neuroscience theories of consciousness.

Specifically, as we have mentioned, the CTM is inspired by cognitive neuroscientist Bernard Baars' theater model of consciousness (Baars, Bernard J., 1997), also called the global workspace (GW) theory of consciousness. However, the CTM is not a standard GW model. The CTM differs from GW in a number of important ways:

- In CTM, competition for global broadcast is formally defined and completely does away with the ill-defined Central Executive of all other GW models;
- In CTM, special processors, including especially its Model-of-the-World processor, construct and employ models of its (*inner* and *outer*) worlds;
- *Brainish*, CTM's multimodal internal language, enables communication between processors and provides the multimodal labeling of sketches in CTM's world models;
- In CTM, *predictive dynamics* (cycles of prediction, testing, feedback and learning, locally and globally), used for machine learning, constantly improves the world models.

The CTM also interacts with its outer world via input *sensors* and output *actuators*. To emphasize CTM's embodied, embedded (in the world), enacted (Thompson, 2007) and extended mind, we use **rCTM** to denote a robot with a CTM brain. The robot could be a human if it has a CTM brain.

While working on this paper, we became aware of Kevin Mitchell's blog post in *Wiring the Brain* (Mitchell, 2023) in which he makes a point similar to one that we make, namely, that many of the major theories of consciousness are compatible and/or complementary. For a similar conclusion, see (Storm, et al., 2024). Even more, Mitchell presents sixteen questions "that a theory of consciousness should be able to encompass". He declares that "even if such a theory can't currently answer all those questions, it should at least provide an *overarching framework*<sup>5</sup> (i.e., what a theory really should be) in which they can be asked in a coherent way, without one question destabilizing what we think we know about the answer to another one."

---

<sup>5</sup> Italics ours.

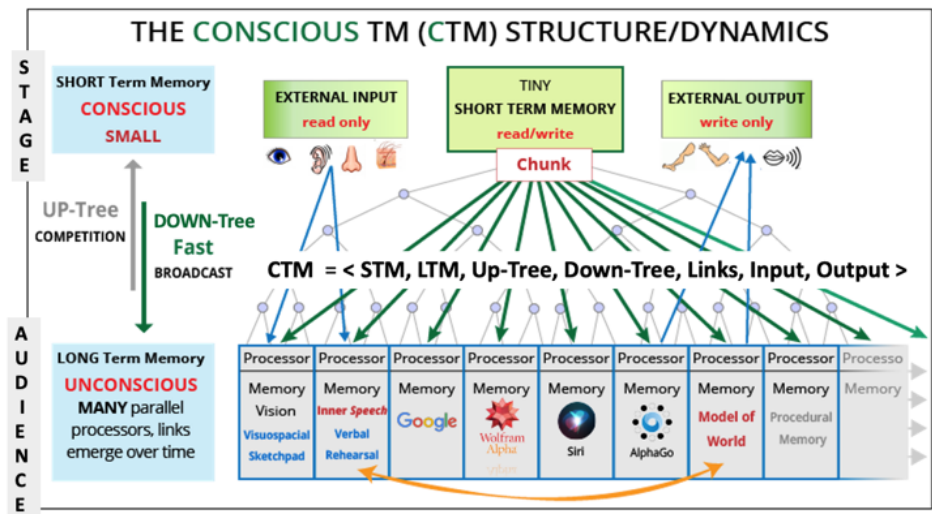
# AI Consciousness is Inevitable

Mitchell’s questions are thoughtful, interesting, and important. Later in this paper (Section 4), we offer preliminary answers from the perspective of rCTM. Our answers to Mitchell’s questions supplement and highlight material in the following brief Overview of rCTM.<sup>6</sup>

## 2 Brief Overview of rCTM, a Robot with a CTM Brain

### 2.1 Formal Definitions of CTM/rCTM

CTM is defined formally as a 7-tuple, (STM, LTM, Up-Tree, Down-Tree, Links, Input, Output). The seven components each have well-defined properties (Blum & Blum, 2022). Here we call it rCTM, a robot with a CTM brain (to emphasize that it is embedded in, and enacts in, its outer world), and outline its properties.



For rCTM, the stage in the theater model is represented by a Short Term Memory (STM) that at any moment in time contains rCTM’s current *conscious content*. STM is not a processor; it is merely a buffer and broadcasting station. The N audience members<sup>7</sup> are represented by a massive collection of (initially independent) powerful processors that comprise rCTM’s principal

<sup>6</sup> In the Overview (Section 2), we annotate paragraphs that refer to Kevin Mitchell’s queries. For example, if a paragraph has a label [KM1], then it refers to Mitchell’s first question, KM1. If Mitchell’s question (in Section 4) is labeled with an asterisk such as KM1\*, then it refers to paragraphs labeled [KM1] in the Overview.

<sup>7</sup> We assume that there are  $N \geq 10^7$  LTM processors (suggested by the  $10^7$  cortical columns in the human brain). Abstractly, these processors are *random access machines*. (*not* Turing machines).

## AI Consciousness is Inevitable

computational machinery and Long Term Memory, together called **LTM**. In the rCTM, all processors are in LTM so when we speak of a processor, we mean an LTM processor.<sup>8</sup> These processors compete to get their information (in the form of chunks, to be defined formally shortly) on stage to be immediately broadcast to the audience.<sup>9</sup> (See Appendix 7.2 for a discussion of rCTM's competition.

rCTM has a *finite* lifetime **T**.<sup>10</sup> At time  $t = 0$ , all but the *Input* and *Output* LTM processors are “generic” with certain basic built-in properties, e.g., some learning/prediction correction algorithms, as well as a preference for choosing the positive over the negative.<sup>11</sup> Their functionalities evolve over time.

But for purposes of building a prototype rCTM, we designate some important LTM processors built in. These include: a *Model-of-the-World processor* (MotWp), actually a collection of processors

---

<sup>8</sup> In rCTM, STM is only a buffer and broadcasting station; all processors are in LTM. Hence, processors like the phonological loop and visuospatial sketchpad (Baddeley & Hitch, 1974), and episodic buffer (Baddeley, 2000) will be in LTM.

<sup>9</sup> As an example of the theater analogy, consider the “What’s her name?” scenario: Suppose at a party, we see someone we know but cannot recall her name. Greatly embarrassed, we rack our brain to remember. An hour later when driving home, her name pops into our head (unfortunately too late). What’s going on?

Racking our brain to remember caused the *urgent* request “What’s her name?” coming from LTM processor **p** to rise to the stage (STM) which immediately broadcasts the question to the audience.

Many (LTM) processors try to answer the query. One such processor recalls we met in a neuroscience class; this information gets to the stage and is broadcast triggering another processor to recall that what’s-her-name is interested in “consciousness”, which is broadcast. Another processor **p’** sends information to the stage asserting that her name likely begins with **S**.

Sometime later the stage receives information from processor **p”** that her name more likely begins with **T** which prompts processor **p””** (who has been paying attention to all the broadcasted information) to claim with *great certainty* that her name is **Tina** – which is correct. The name is broadcast from the stage, our audience of processors receives it, and we finally consciously remember her name. Our conscious self has no idea how her name was found.

(Based on the correct outcome, learning algorithms internal to each processor, cause processor **p’** to lower the importance ( $|\text{weight}|$ ) it gives its information and cause **p”** to increase the importance.)

<sup>10</sup> In rCTM, parameters **T** and **N** (with **T** = lifetime in ticks and **N** = number of LTM processors) are equal: **T = N**. More generally,  $T = cN$  for some small multiplicative constant like  $c = 3$ .

Time  $t = 0, 1, 2, 3, \dots, T$  is measured in discrete clock ticks (The 10 to 100 ticks per second is roughly equal to the alpha and beta/gamma brain wave frequencies).

<sup>11</sup> This built-in preference creates a predilection for survival. Primal chunks of pain and pleasure have negative and positive weights, respectively. Signals from built in nociceptor fibers have built in negative weights. The positive valence comes largely from the relief of pain. See section 2.3.1.

## AI Consciousness is Inevitable

(including all input and output processors) that collaborate with all processors, to build models of rCTM's inner and outer worlds<sup>12</sup>; *Sensory* processors (with **input** from rCTM's outer world via its various *sensors*<sup>13</sup>); *Motor* processors (with **output** to rCTM's outer world via motor *actuators*<sup>14</sup>); and so on.

We also allow off-the-shelf processors (like ChatGPT, Google and WolframAlpha) that *extend* rCTM's capabilities at the start.

While each processor may have its own distinct language, processors communicate with each other in *Brainish*, Brainish being rCTM's multimodal inner language (words and grammar). Brainish *gists*<sup>15</sup> are succinct Brainish words and phrases that fuse modalities (e.g., sight, sounds, smells, tactile) and processes.

A Brainish gist is like a frame of a dream. The Brainish language evolves over rCTM's lifetime, from nothing initially into an ever growing dictionary of words and ever improving simple creole-like grammar Brainish can differ from one rCTM to another. In section 2.3.1 we indicate how Brainish evolves.<sup>16</sup>

*Input/Sensory* processors convert input information into Brainish gists. *Output/actuator* processors convert Brainish gists into commands for rCTM's actuators.

LTM processors *compete* in a well-defined (*fast* and natural) probabilistic *competition* (Appendix 7.2) to get their questions, answers, and information in the form of a *chunk* onto the stage (STM).

---

<sup>12</sup> rCTM's inner world is what's inside rCTM itself, including rCTM's processors and its "thoughts" and memories; rCTM's outer world is its environment, which includes rCTM itself.

<sup>13</sup> Ears, eyes, nose, skin, mouth/taste, ... .

<sup>14</sup> Arms, hands, legs, mouth/vocal actuators, ... .

<sup>15</sup> We were unaware of Jeremy Wolfe's 1998 paper (Wolfe, 1998) until we read Fei-Fei Li's wonderful book, *The Worlds I See* (Li, 2023), where she identifies Wolfe as the "gist" guy. Wolfe uses the word "gist" in reference to "what we remember – and what we forget – when we recall a scene." Our explanation of *change blindness* in the CTM (Blum & Blum, 2022) is almost identical to Wolfe's explanation of change blindness in humans: the same gist describes both the original and changed scenes.

<sup>16</sup> Paul Liang is developing a computational framework for Brainish based on multimodal machine learning (Liang, 2022). This is different but aligned with rCTM's evolved theoretical Brainish. See section 2.3.1.

## AI Consciousness is Inevitable

The competition is hosted by the **Up-Tree**, a *perfect*<sup>17</sup> binary tree of height **h** which has a leaf in each LTM processor and its root in STM. At each clock tick, *a new competition starts* with each processor putting a *chunk* of information into its Up-Tree leaf node.

A *chunk* is defined formally to be a tuple, **<address, time, gist, weight, auxiliary information>**, consisting of (in order of importance): a succinct Brainish gist of information; a *valenced* weight (to indicate the importance/urgency/value/confidence/sentiment (+/-) the originating processor assigns the gist); the address<sup>18</sup> of the originating processor; the time the chunk was created; plus some auxiliary information. (See Appendix 7.2 for a discussion of auxiliary information.)

Each submitted chunk competes locally with its neighbor (as in a tennis or chess tournament). In a clock tick, the local winning chunk is chosen and moves up one level of the Up-Tree, ready to compete with its neighbor. At each clock tick, there are active local competitions at every node in every level from **0** to **h-1** of the Up-Tree. The competition that begins at time **t** ends at time **t+h** with the *winning chunk* in STM.<sup>19</sup>

Notably, for the probabilistic rCTM competition, it is proved that a chunk wins the competition with probability proportional to (a *function* of) its weight. As a consequence, the winning chunk is independent of the processor's location! Clearly this is an important feature for a machine or brain; no moving around of processors is needed. (It is a property that would be difficult if not impossible to achieve in fair tennis or chess tournaments where original pairings depend on current rankings.) (See Appendix 7.2 for a fuller discussion of the Up-Tree competition.)

The chunk that gets onto the stage (STM), i.e., the *winning chunk*, is called rCTM's current *conscious content* and is immediately *globally broadcast* (in one clock tick) via the **Down-Tree** (a bush of height **1** with a root in STM and **N** branches, one leaf in each LTM processor) to the audience (of all LTM processors). [KM2] [KM5]

---

<sup>17</sup> A *perfect* binary tree is a binary tree in which all leaf nodes are at the same depth. This depth is also the height of the tree. If **h** is the height of a perfect binary tree, then the tree has **N = 2<sup>h</sup>** leaves. Each node, except the root node, has a unique sibling neighbor. For simplicity, we choose a perfect binary tree.

<sup>18</sup> If there are  $\sim 10^7$  processors, the address is a  $\sim 7$  digit number.

<sup>19</sup> If the interval between clock ticks is  $10\text{ms}$  and the number of LTM processors is  $\sim 10^7$ , then a full competition takes about  $230\text{ms}$ .



## AI Consciousness is Inevitable

The single chunk in STM to be globally broadcast will enable rCTM to focus attention on the winning gist. “One” is not the “magical number”  $7 \pm 2$ , but we are looking for simplicity and a single chunk will do.<sup>20</sup>

### 2.2 Conscious Attention in rCTM

**Formal definition 1.** *Conscious attention* in rCTM is the *reception* by all LTM processors of the broadcast of rCTM’s current conscious content. [KM2] [KM5]

In other words, rCTM *pays conscious attention* at time  $t+h+1$  to the winner of the competition that commenced at time  $t$ .<sup>21</sup>

**We call a sequence of receptions of broadcasted chunks, a *stream of consciousness*.**

The rCTM has no links at birth. A 2-way **link** forms between processors A and B when A acknowledges to B that B has broadcast useful information for A.<sup>22</sup> Such links enable *conscious communication*, i.e. communication that goes through STM, to be replaced by more direct and faster *unconscious communication* through links. Thus, when rCTM initially learns to ride a bike, most communication is done consciously until relevant processor links have formed. Then, for the most part, riding a bike is done unconsciously until an obstacle is encountered, forcing rCTM to pay conscious attention again. [KM6]

Built into each LTM processor, in particular the MotWp, are algorithms that work unconsciously to make predictions and correct itself. The MotWp plays an important role in planning, predicting, exploring, testing and correcting/learning. More generally, all LTM processors make *predictions* and get *feedback* from broadcasts, from linked processors, and from the outer world. Based on this

---

<sup>20</sup> The “magical number”  $7 \pm 2$  was proposed by George Miller (Miller G. A., 1956) as the number of “chunks” that a human can hold in their “short term memory” at any moment of time.

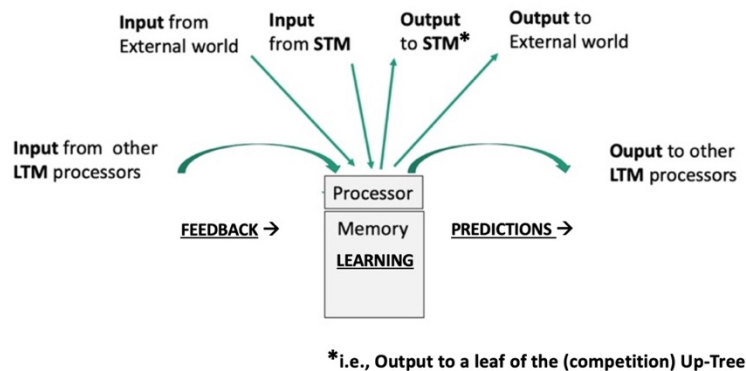
<sup>21</sup> There is a delay of  $h$  clock ticks between when {LTM processors enter their chunks into the competition for STM} and {the winner is revealed as rCTM’s “conscious content”}. In one more clock tick, rCTM pays conscious attention to the winner. This  $h+1$  delay is somewhat analogous to behavioral and brain studies going back to (Libet, 1985) that suggest a delay between when unconscious processors in our brains make decisions and when we become conscious of them.

<sup>22</sup> The rCTM has no links at birth. A 2-way link forms between processors A and B when A acknowledges to B that B has broadcast useful information for A. In our earlier “What’s her name?” scenario (footnote in section 2.1), when processor  $p$  sees that processors  $p''$  and  $p'''$  have useful information for it, and  $p$  acknowledges their usefulness,  $p$  forms bi-directional links with  $p''$  and with  $p'''$ .

## AI Consciousness is Inevitable

feedback, *learning* algorithms internal to each processor use prediction errors to correct and improve that processor's behavior.

A major rCTM goal is to ensure that its predictions (i.e., its processor's predictions) are kept as accurate as possible. For the most part this is done unconsciously.



Learning algorithms include each processor's built-in *Sleeping Experts Learning* algorithm (necessarily modified for rCTM) that adjusts the weights each processor gives its gists. See (Blum A. , 1995) and (Blum, Hopcroft, & Kannan, 2015). [KM7]

Thus, we can already see some basic *predictive dynamics* (prediction + testing + feedback + learning/correction) occurring locally and globally within rCTM. [KM2]

rCTM's competition, broadcast, attention and immediate direct communication via links is reminiscent of a process that Dehaene and Jean-Pierre Changeux call *ignition* (Dehaene & Changeux, 2005).

But... for *feelings of consciousness*, attention is not all you need. More is required.

### 2.3 Conscious Awareness and the Feeling of Consciousness in rCTM

How might rCTM, a formal machine model, experience feelings such as pleasure and pain? More generally, how might rCTM experience *feelings of consciousness* and *conscious awareness*?

In this section we look at the central role played by rCTM's *Model-of-the-World processor* (MotWp) in the dynamic evolution of rCTM's subjective consciousness.

The MotWp, *in collaboration with* all processors, creates *models* of rCTM's worlds, inner and outer, collectively called the *Model-of-the-World* (MotW). The MotWp, a collection of processors,

## AI Consciousness is Inevitable

includes rCTM's Sensory and Motor processors in order for it (the MotWp) to receive inputs directly from the outer world via sensors and send outputs directly to the outer world via actuators.

The MotW represents rCTM's current and continuing view of rCTM's worlds, inner and outer. This information becomes conscious when chunks containing gists about the MotW win the competition for STM and are received by all processors. The world that rCTM consciously "sees" or more generally "senses" or "knows" *is* the MotW, as **rCTM is not conscious of anything that does not come directly from STM**. This will be especially important for understanding the *feeling of consciousness* in rCTM! [KM1] [KM3] [KM16]

Specifically, the MotWp constructs *sketches* in the MotW, succinct Brainish descriptions (representations) of *referents* in rCTM's worlds (inner and outer). A referent might be a red rose in rCTM's outer world or a feeling of pleasure in its inner world or a thought. Their sketches may be *labeled* with succinct *Brainish gists* such as: BRIGHT\_COLOR/SWEET\_SMELL/SILKY\_TOUCH or FEELING\_HAPPY. These *labels* are mnemonic indicators of what rCTM "learns", "senses" or "feels" about those referents.<sup>23</sup> The slashes are mnemonic indicators of fused modalities.<sup>24</sup>

Sketches and their labels evolve over time. In particular, the sketch "rCTM" in the MotW whose referent is rCTM itself<sup>25</sup> will develop from scratch and eventually be labeled with SELF/CONSCIOUS/ FEELS, etc.. [KM1] [KM3] [KM9]

We now look at how the MotW and Brainish co-evolve to produce rCTM's evolving feelings of subjective consciousness.

---

<sup>23</sup> These labels are mnemonic indicators of what we call *qualia* in rCTM, the subjective character of rCTM's conscious experiences.

<sup>24</sup> Formally, fusion will be represented by a collection of pointers to previously defined multimodal words and phrases.

<sup>25</sup> We are often asked, isn't this process recursive? Doesn't the sketch of rCTM have a sketch of rCTM have a sketch of rCTM, etc. ? Yes, up to a point. But, at each iteration the current sketch is degraded, so the process rapidly becomes null.

## AI Consciousness is Inevitable

### 2.3.1 The Model-of-the-World and Brainish co-Evolve

Both the Model-of-the-World and Brainish evolve over time and play an essential role in the feeling of “what it is like” to be an rCTM.<sup>26</sup> [KM1] [KM3] [KM9]

The infant rCTM (which we think of as a simulation of a human infant with a CTM brain) has only a grainy foggy MotW, which does not even include a sketch of itself. Sketches develop over time, become refined, and get Brainish labels.

For example, at some point in time, the MotW will contain a rough sketch of the infant’s left leg. When the infant rCTM discovers it can move its left leg (an actuator) by the *power of thought*, the MotWp appends the label SELF to the sketch of the left leg.<sup>27</sup> In this way, the foggy MotW becomes more accurate. [KM3] [KM13][KM16]

To perform an action by the **power of thought** is to perform the action in the MotW (as in a dream), then confirm that it actually got done in the world. In this case, the MotWp moves the left leg sketch in its MotW (this simulation is a kind of prediction), which sends a command to the left leg actuator to move. The MotWp detects that movement via rCTM’s sensors and, by repeating the action, becomes more certain that rCTM is itself responsible for moving the leg (that is, becomes more certain that its prediction is correct). Now convinced, it labels the left leg sketch with SELF. [KM13]

Feelings of pain and pleasure in rCTM start to develop even earlier. They are key exemplars of the Hard Problem. We indicate how these feelings come about in rCTM:

Consider an infant rCTM at the moment it is born. A processor that monitors the O<sub>2</sub> level, the O<sub>2</sub>\_Gauge processor, raises its growing concern for the lack of O<sub>2</sub> by submitting a sequence of chunks having negatively valenced weights of increasingly high absolute value to the Up-Tree competition. At some point, that O<sub>2</sub>\_Gauge processor’s chunks get onto the stage (STM), their huge

---

<sup>26</sup> Thomas Nagel’s “what it is like” (Nagel, 1974) is often taken to be the canonical definition of phenomenal consciousness.

<sup>27</sup> Certain pathologies will occur if a breakdown in rCTM causes its MotWp to mislabel. For example, if the sketch of that leg gets labeled NOT-SELF at some point, rCTM might beg to get its leg amputated, even if it still functions properly. This would be an example of body integrity dysphoria (body integrity identity disorder) in rCTM. Other pathologies due to faulty labeling in the MotW include: phantom limb syndrome (a sketch of an amputated arm actuator is mislabeled SELF), Cotard’s syndrome (the sketch of rCTM is labeled DEAD), paranoia (the sketch of rCTM’s best friend is labeled SPY), ... [K10]

## AI Consciousness is Inevitable

and growing negative weight signals a desperate “scream” for something ( $O_2$ ). All processors hear this “scream”.

Every processor is programmed to store the  $O_2$ \_Gauge processor’s address  $\mathbf{p}$  and current weight  $\mathbf{w}$  as a Brainish word  $(\mathbf{p}, \mathbf{w})$  when it receives a broadcasted chunk from the  $O_2$ \_Gauge processor. Every processor is also programmed, in case  $\mathbf{w}$  has a negative valence, to give fullest possible attention to reducing a huge  $|\mathbf{w}|$ .

For such  $\mathbf{w}$ , we interpret the Brainish word  $(\mathbf{p}, \mathbf{w})$  as FEELS\_PAIN\_FROM\_THE\_LACK\_OF\_ $O_2$  which will become generic for FEELS\_PAIN, rCTM’s first Brainish word. The MotWp *labels* the *sketch* of the  $O_2$ \_Gauge processor with the *current pair*  $(\mathbf{p}, \mathbf{w})$ .

The processors have absolutely no idea what to do (as the infant was just born) to reduce that huge  $|\text{weight}|$ . They know only that they must do something. Processors that control rCTM’s actuators command them to do something, anything. The arms and legs flail. The infant pees and poos. The Motor/Vocal processor commands the vocal actuator to scream and cry. This last one works! The cry opens the “lungs”, the infant takes its first breath, the weights in chunks generated by the  $O_2$ \_Gauge processor then return to normal.

The next time the infant rCTM needs help, screaming and crying works again. In short order, the infant rCTM learns that a good response to pain of any kind is to scream and cry.

Recapping, we have an early example of how the feeling of pain, the Brainish word for pain, and rCTM’s consequent response, arise in rCTM:

All processors pay *conscious attention* to the  $O_2$ \_Gauge processor’s desperate call for  $O_2$ . This *triggers* (ignites) the MotWp to append the  $O_2$ \_Gauge processor’s (address  $\mathbf{p}$ , current chunk weight  $\mathbf{w}$ ) which we interpret as FEELS\_PAIN, to a sketch of the  $O_2$ \_Gauge processor in the MotW. This also *triggers* (ignites) all processors to do something, and they do anything they know how to do. The Vocal processor/vocal actuator team solves this primal problem with a screaming cry, which the rCTM also learns to be a good response to any kind of pain. The Brainish word  $(\mathbf{p}, \mathbf{w})$  that triggered the primal cry is rCTM’s first word for pain.

## AI Consciousness is Inevitable

This dynamic interaction between conscious attention, the Brainish labeling of sketches in the MotW and the triggering of processors to action (ignition) is an example of how rCTM experiences the subjective feeling of pain.<sup>28</sup>

When the infant rCTM becomes a toddler and skins its knee, its MotW will hold a sketch of a bloody knee labeled SELF/FEELS\_PAIN. In a disorder in which that knee is labeled SELF but not PAIN, or PAIN but not SELF, the rCTM has *pain asymbolia*: it knows there is pain but does not suffer from it.<sup>29</sup>[KM14]

Let's return to the infant rCTM and suppose the infant rCTM's Fuel Gauge is low. The Fuel Gauge processor creates a chunk with |weight| proportional to the need. When this chunk wins the competition and gets globally broadcast, the Procure Fuel processor activates (getting the vocal actuator to cry out) causing the Fuel Source (mother) in the outside world to respond. The infant rCTM learns that the Fuel Source relieves FEELS\_PAIN when the Fuel Gauge indicates low fuel.

The current sketch of the Fuel Gage processor in the MotW is labeled FEELS\_PAIN/HUNGER. The Fuel Source sketch (of mother) in the MotW is labeled RELIEVES PAIN/RELIEVES HUNGER, and will also get labeled PLEASURE\_PROVIDER. When a chunk with a gist containing the Fuel Source sketch and its label PLEASURE\_PROVIDER/RELIEVES\_PAIN and a positively valenced weight is broadcast, the infant rCTM learns that PLEASURE relieves PAIN, and this reinforces its preference for positively valenced gists over negative ones.<sup>30</sup> [KM3][KM12]

Over time and with experience, these primitive and primal Brainish words, labels and sketches will change and become richer. The early assigned valences with preference for the positive will apply to an increasing number of situations.

---

<sup>28</sup> We note that the flailing *process* itself contributes to rCTM's feeling of pain (and desperation). A mnemonic for this process is fused with the primal word for pain.

<sup>29</sup> A person who has pain and *knows* everything about it but lacks the ability to *feel* its agony has *pain asymbolia*. Such a person is not motivated to respond normally to pain. Children born with pain asymbolia rarely live past the age of 3. The experience of pain, whether physical or emotional, serves as a motivator for responding appropriately to the pain. See (Grahek, 2001; 2007) and (Philip Gerrans, 2024).

<sup>30</sup> When a human mother gives a breast to her infant, the infant learns that the breast relieves the pain of hunger. The breast gets incorporated into the infant's MotW labeled with RELIEVES\_PAIN/FEELS PLEASURE. Unless brought up in a psychotic household, the human infant learns that pleasure relieves pain and prefers pleasure over pain. This is one of many ways in which pain and pleasure are not symmetrical.

## AI Consciousness is Inevitable

As Brainish words evolve, the Brainish language will develop a simple basic (creole-like) grammar<sup>31</sup>. [KM1]

### 2.3.2 Conscious Awareness in rCTM

*Conscious awareness* in rCTM is the dynamic interaction between conscious attention and an evolving MotW with increasingly rich Brainish-labeled sketches.<sup>32</sup>

**Formal definition 2.** rCTM becomes *consciously aware* of a MotW Brainish-labeled sketch (sketch and label) when it pays conscious attention to a chunk containing a gist with this labeled sketch.

Gists in these chunks are like frames in a dream.<sup>33</sup> [KM2]

For example, rCTM becomes consciously aware of a red rose in the outer world when it pays conscious attention to the sketch of the rose in the MotW with its label BRIGHT\_RED/SMELLS\_SWEET/FEELS\_SILKY.

As rCTM becomes increasingly consciously aware, some processor will call attention to this fact and many others will concur. The MotWp will label the sketch of rCTM in the MotW as CONSCIOUSLY\_AWARE/FEELS\_CONSCIOUS or simply CONSCIOUS.<sup>34</sup>

Looking at rCTM from the viewpoint of the outside world, we see that something about rCTM is conscious; specifically, the rCTM considers itself conscious. What is conscious cannot be the MotWp or any other processor, as processors have no feelings; they are just machines running algorithms. Our proposal that rCTM as a whole “feels it is conscious” is a consequence in part of

---

<sup>31</sup> For information on simple (creole-like) languages see: (McWhorter, 1998), (McWhorter, 2008), (Sigal, 2022) and (Bancu, et al., 2024).

<sup>32</sup> That is, conscious awareness is the interplay between focused attention on specific aspects of the entity’s inner and outer worlds, and the continuous updating and interpretation of this information within the entity’s world models that contain multimodal linguistic labels of referent sketches.

<sup>33</sup> rCTM *dreams* are streams of consciousness generated when the input sensors and output actuators are inactive, and rCTM’s *Dream* processor gets to work. Although dreams are “felt” as real, they can also be fantastical since their predictions are not being tested in the world. We propose that (testing for) *dreaming* is a (partial) *test for subjective consciousness*.

<sup>34</sup> What makes CONSCIOUS the Brainish word for “conscious” is that the address-weight pair of the processor that called attention to rCTM’s conscious awareness is CONSCIOUS in Brainish.

## AI Consciousness is Inevitable

the fact that the MotWp views the “rCTM” in its MotW as conscious, and that this view is broadcast to all processors. This *is* rCTM’s *subjective consciousness*.

### 2.4 rCTM as a Framework for Artificial General Intelligence (AGI)

Before indicating rCTM’s alignment with a number of major theories of consciousness, we remark on rCTM’s potential to serve as a framework for constructing an Artificial General Intelligence (AGI). This is a result of rCTM’s global architecture (kindred to arguments made for global latent workspace by (VanRullen & Kanai, 2021)) and, at the same time, the result of an essential difference between rCTM and Baars’ global workspace. rCTM has *no* Central Executive. This is a feature, not a bug. It enables rCTM to become an AGI (Blum & Blum, 2023):

The competition to get information on stage considers the |weight|ed information submitted by a huge collection of ( $N$ ) processors. And it does this quickly ( $\log_2 N$  steps). This enables rCTM to engage processors to solve problems, even though rCTM does not know which of its processors might have the interest, expertise, or time to do so. [KM15]

Specifically, if rCTM, meaning one (or more) of its LTM processors, has a problem to solve, the processor can submit the problem to the competition as a chunk with high enough |weight| giving it a high probability of winning and thus being globally broadcast to all processors. Processors with the interest, expertise and time to work on the problem will respond with appropriately |weight|ed chunks. In this way, ideas from unexpected sources may contribute to solving the problem, and useful collaborations can emerge.<sup>35</sup> [KM15]

A Central Executive would have to know which processors had the inclination, expertise, and resources to solve problems as they arise, and figure this out quickly. Baars’ does not say how a Central Executive could do this. [KM15]

More generally, we predict that a Central Executive is not needed for consciousness or for general intelligence; indeed, it might be an impediment.

To our knowledge, all other models of consciousness have a Central Executive.

---

<sup>35</sup> If rCTM’s competition is based on the simple  $f$  value = intensity (see Appendix 7.2), the processor can submit the problem to the competition with |weight| = intensity of the currently broadcasted chunk, which is the sum of all |weights| of all submitted chunks. In that case, the chunk with that high |weight| will win, with probability  $\geq 1/2$  and become globally broadcast to all processors.



### 3 Alignment of rCTM with Other Theories of Consciousness

We have presented an overview of the rCTM model. Now we indicate how, at a high level, the model naturally aligns with and integrates key features of major theories of consciousness, further supporting also our view that rCTM provides a framework for building a conscious machine.

#### 3.1 Global Workspace (GW)/Global Neuronal Workspace (GNW)

rCTM aligns broadly with the architectural and global broadcasting features of the *global workspace* theory of consciousness (Baars, Bernard J., 1997), and at a high level with the *global neuronal workspace* theory of consciousness of neuroscientists Stanislas Dehaene, Jean-Pierre Changeux (Dehaene & Changeux, 2005), (Dehaene S., 2014), and others.<sup>36</sup>

However, rCTM differs from GW in significant ways. For example: rCTM has a formally defined natural competition for information to become globally broadcast; it constructs world models; and rCTM has no Central Executive, a feature, not a bug.

#### 3.2 Attention Schema Theory (AST)

rCTM's ability to construct and utilize models of rCTM's worlds (inner and outer), and the key role they play in rCTM's *conscious awareness*, align closely with neuroscientist Michael Graziano's *attention schema theory* of consciousness (Graziano, Guterstam, Bio, & Wilterson, 2020). AST propose that the brain is an information processing machine that constructs a simplified model of attention, just as it constructs a simplified model of the body, the Body Schema. According to AST, this Attention Schema provides a sufficiently adequate description of what it is attending to for the brain to conclude that it is "aware".<sup>37</sup>

#### 3.3 Predictive Processing (PP)

Predictive processing asserts that the brain is constantly inferring, correcting and updating its predictions, generally based on motor outputs and sensory inputs. rCTM's *predictive dynamics*

---

<sup>36</sup> Additional references for GNW include: (Dehaene & Naccache, 2001), (Sergent & Dehaene, 2005), (Dehaene & Changeux, 2011), and (Mashour, Roelfsema, Changeux, & Dehaene, 2020).

<sup>37</sup> In this regard, rCTM and AST both align with philosophers Daniel Dennett's and Keith Frankish's view that phenomenal consciousness is "an illusion" if understood by their frequent explications: "Consciousness is the brain's interface for itself." (Dennett); and "Phenomenal consciousness is real, but not what you think it is." (Frankish).

## AI Consciousness is Inevitable

(cycles of prediction, testing, feedback, and learning/ correcting), locally and globally, align with various incarnations of *predictive processing* (von Helmholtz, 1866; 1962), (Friston K. , 2010), (Cleeremans, 2014), (Clark A. , 2015), (Hohwy & Seth, 2020), and others.<sup>38</sup>

### 3.4 Embodied Embedded Enactive Extended Mind (EEEE Mind)

rCTM's ability to construct (and utilize) models of its worlds containing rich Brainish-labeled sketches (leading to its *feelings of consciousness*) derives in part from its embodied, embedded, enactive, and extended (EEEE) mind.

This aligns with the “4E” view that consciousness, like cognition (Carney, 2020), involves more than brain function (Rowlands, 2010). For consciousness, the 4E's are:

- *Embodied*: Incorporating relations with the entity's *body parts* and *processes* is essential for phenomenal consciousness. See, (Damasio, 1994), (Edelman, 2006) and (Shanahan, 2005).<sup>39</sup>
- *Embedded*, and *Enactive*: Being *embedded in the outer world* and *enacting*/interreacting with it, thus affecting the world and creating experiences, is necessary for phenomenal consciousness. See, (Maturana & Varela, 1972), (Maturana & Varela, 1980), (Varela, Thompson, & Rosch, 1991), (Thompson, 2007), and (Clark A. , 2008).
- *Extended*: Consciousness is further enhanced by the entity having access to considerable external resources (such as libraries, Google, ChatGPT, Mathematica, ...). See, (Clark & Chalmers, 1998).

rCTM is *embedded* in its outer world and, through its *embodied* actuators, can *enact* in this world, thus influencing what it senses and experiences. rCTM's “mind” is *extended* by information it gets from resources in its outer world, and from its embedded (or linked) off-the-shelf processors.

---

<sup>38</sup> Other references include: (McClelland & Rumelhart, 1981), (Lee & Mumford, 2003), (Friston K. , 2005), (Clark A. , 2015), (Seth, 2015), (Miller, Clark, & Schlicht, 2022).

<sup>39</sup> We note that here Shanahan views the global workspace as key to access consciousness, but that phenomenal consciousness requires, in addition, embodiment.

### 3.5 Integrated Information Theory (IIT)

IIT, the theory of consciousness developed by Giulio Tononi (Tononi, 2004), and supported by Koch (Tononi & Koch, 2015), proposes a measure of consciousness called Phi that, in essence, measures the amount of feedback and interconnectedness in a system. rCTM’s extensive feedback (its predictive dynamics, globally and locally) and its interconnectedness (global broadcasts and its fused multi-modal Brainish gists) contributes to a high Phi.

### 3.6 Evolutionary Theories of Consciousness

rCTM aligns with aspects of *evolutionary theories* of consciousness.

Oryan Zacks and Eva Jablonka provide evidence for the evolutionary development of a modified *global neuronal workspace* in vertebrates (Zacks & Jablonka, 2023) reinforcing our suggestion that an AI with a global workspace architecture could possess access consciousness.<sup>40</sup>

In *Sentience*, Nicholas Humphrey presents an evolutionary argument for the development of phenomenal consciousness in warm-blooded animals (Humphrey, 2023). In “The Road Taken” (Chapter 12 of *Sentience*), Humphrey spins a “developing narrative”, starting with “a primitive amoeba-like animal floating in the ancient seas. Stuff happens. ...” The resulting story provides a roadmap for how an entity might create world models and a sense of self. Indeed, this roadmap closely parallels the way rCTM’s world models evolve, and how rCTM develops its sense of self and subjective conscious awareness.<sup>41</sup> See section 2.3.12.3.1 where we indicate how Brainish and the Model of the World co-evolve.

Thus, while the former evolutionary theory (Zacks & Jablonka, 2023) aligns with rCTM’s built-in architecture, the latter theory (Humphrey, 2023) aligns with rCTM’s development of subjective consciousness over time.

---

<sup>40</sup> See (Ginsburg & Jablonka, 2019) for an extensive treatise on the evolutionary development of consciousness and their Unlimited Associative Learning (UAL) theory of consciousness.

<sup>41</sup> We claim Humphrey actually gives a road map for how an entity, warm blooded or not, might create world models and sense of self. As an exercise, we have re-written part of Chapter 12 (Humphrey, 2023), “The Road Taken”, from the perspective of rCTM and sent a copy to Humphrey. His reply, “It would be great if we could meld these theories.” (Personal communication with Nick Humphrey, Oct 9, 2023.)

### 3.7 Extended Reticulothalamic Activating System (ERTAS) + Free Energy Principle (FEP)

In *The Hidden Spring*, Marc Solms makes the case that the source of consciousness is the arousal processes in the upper brain stem (Solms M. , 2021). More generally, Solms cites the Extended Reticulothalamic Activating System (ERTAS) as the generator of feelings and affects, enabling consciousness. “Affective qualia” is the result of homeostasis. “Deviation away from a homeostatic settling point (increasing uncertainty) is felt as unpleasure, and returning toward it (decreasing uncertainty) is felt as pleasure” (Solms M. , 2019). Homeostasis arises by a system resisting entropy, i.e., minimizing free energy (Solms & Friston, 2018). This is enabled by a *Markov blanket* (containing the system’s input sensors and output actuators) that insulates the internal system from its outer world. The “system must incorporate a *model of the world*, which then becomes *the basis upon which it acts*<sup>42</sup>” (Solms M. , 2019). Similarly, rCTM is consciously aware only of its MotW broadcasted chunks, and thus the MotW becomes “the basis upon which it acts”.

Kevin Mitchell points out (personal communication) that another important point from Solms is that the ascending homeostatic signals, which track different needs, must be valenced but also must have some distinguishing “qualities” so that when they are submitted to central decision-making units, the sources of the signals can be kept track of - so the organism doesn't mistake feeling thirsty for feeling tired (both of which feel BAD)”. In rCTM, gists are the distinguishing qualities, and chunks contain addresses of originating processors, so tracking is implicit. No need for a central decision-maker.

At a high level, rCTM aligns with ERTAS + FEP:

- Although GW models generally consider processors as performing cortical functions, rCTM goes beyond that. There is nothing to preclude rCTM from having processors that function as the ERTAS.
- In (Blum & Blum, 2021), we discuss pleasure and pain in the CTM (known here as rCTM). Our discussion of pleasure aligns with Solms, and also with (Berridge & Kringelbach, 2015). In (Blum & Blum, 2021) we discuss in more detail how feelings of pain might be generated.

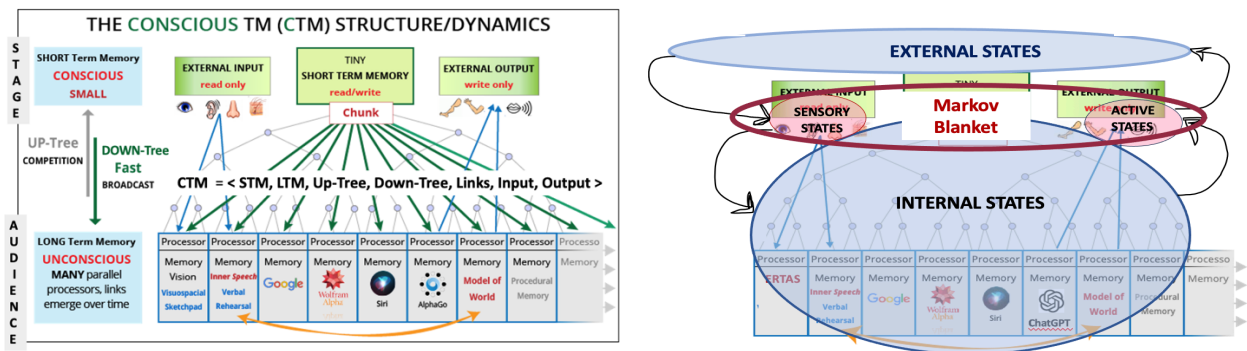
---

<sup>42</sup> Italics here ours.

## AI Consciousness is Inevitable

- Predictive dynamics (cycles of prediction, testing, feedback and correcting/learning) in rCTM works to reduce prediction errors, an analogue to minimizing free energy.
- And the incorporated Model-of-the-World in rCTM is *the basis upon which rCTM acts*.
- Tracking of similar feelings is naturally maintained by chunks that contain gists, weights, valences, address of originating processor, and time created.

Evocatively, the well-known Friston diagram (Parr, Da Costa, & Friston, 2019), with Markov blanket separating internal and external states, is clearly realized in rCTM:



## 4 Addressing Kevin Mitchell's questions from the perspective of rCTM

Here we address Kevin Mitchell's sixteen questions (Mitchell, 2023)<sup>43</sup> from the perspective of the robot with the CTM brain (rCTM).

Our answers refer to and supplement what we have discussed in our Overview (Section 2). They deal *only* with the rCTM model, meaning an entity with a CTM brain. These answers say nothing about other models. They say nothing about whether a worm is conscious or not - unless the worm has a CTM brain. From here on, unless otherwise stated, everything we have to say is about the rCTM model.




<sup>43</sup> Many of Mitchell's questions are in fact a collection of intertwined questions.

## AI Consciousness is Inevitable

In the following, Mitchell's questions are printed in **bold**. Our answers follow in non-bold print. In the Overview (Section 2), we annotated paragraphs that refer to Kevin Mitchell's queries. For example, a paragraph labeled [KM1] refers to Mitchell's first question, KM1.<sup>44</sup>

**KM1\*. What kinds of things are sentient? What kinds of things is it like something to be? What is the basis of subjective experience and what kinds of things have it?**

rCTM is *sentient*, meaning it is able to sense and feel things. As mentioned above, we have nothing to say about entities that are not rCTMs. However, we can and will sometimes say (as we do below) what parts of rCTM are responsible for its sentience.

The ability of rCTM to construct models of its worlds (inner and outer) is fundamental for its *subjective experiences*. These experiences are described in the model not with English words but with multimodal *Brainish-labeled sketches* of referents in rCTM's worlds. In the case of a red rose, the label, a Brainish word, might be an image  fused with its odors , smooth touch , and such. We suggest that Brainish grammar would be like a simple creole language.

The Model-of-the-World processor (MotWp) and the Model-of-the-World (MotW)<sup>45</sup> it creates play an essential role in "what it is like" to be a rCTM. The sketches and labels (as well as Brainish itself) are *created and evolve* throughout the life of rCTM. The labels succinctly indicate what rCTM learns, senses or feels about the referents. For example, the label SELF applied to a sketch in the MotW indicates that that particular sketch's referent is felt as (a part or the whole of) rCTM itSELF.

**KM2\*. Does being sentient necessarily involve conscious awareness? Does awareness (of anything) necessarily entail self-awareness?. What is required for 'the lights to be on'?**

In this paper we define two related notions of consciousness in rCTM, conscious attention and conscious awareness. We review these formal rCTM definitions here:

*Conscious attention* (access consciousness) in rCTM occurs when all LTM processors receive the global broadcast of rCTM's current *conscious content*, that being the current winning chunk in the competition for STM.

---

<sup>44</sup> If Mitchell's question is labeled with an asterisk such as KM1\*, then it refers to paragraphs labeled [KM1] in the Overview.

<sup>45</sup> The **MotWp** is a collection of processors. Likewise, the **MotW** (no **p**) is a collection of (inner and outer) world models.

## AI Consciousness is Inevitable

*Conscious awareness* (subjective consciousness) in rCTM is the dynamic interaction between conscious attention and an evolving MotW that has Brainish. (Conscious awareness arises in rCTM when the broadcasted chunk refers to a Brainish-labeled sketch in the MotW. The labels describe what rCTM is consciously aware of.)

In rCTM, the notion of sentience aligns more with our definition of conscious awareness than with conscious attention.

**Does awareness necessarily entail *self-awareness*?** No. The infant rCTM initially builds a world model that does not include a labeled sketch of itself, so it has no self-awareness. In time, however, that model will include a rough labeled sketch of itself and the label SELF. The label SELF marks the beginning of self-awareness, which eventually develops into full-blown self-awareness.

In rCTM, the *lights come on* gradually, as the MotW gets populated with sketches and their labels. (For more on this, see our answer to KM4.)

**KM3\*. What distinguishes conscious from non-conscious entities? (That is, why do some entities have the *capacity* for consciousness while other kinds of things do not?) Are there entities with different degrees or kinds of consciousness or a sharp boundary?**

To respond to these questions, we replace “entities” with “rCTM”.

rCTM pays conscious attention to every broadcast. However, *absent a* Model-of-the-World-processor (MotWp) and its Model-of-the-World (MotW) with sketches labeled in Brainish, there is *no* conscious awareness. As stated above, *conscious awareness* arises when the broadcasted chunk refers to a Brainish-labeled sketch in the MotW. The labels describe what rCTM is consciously aware of.

rCTM can be conscious (in both senses) when awake or dreaming. It is not conscious when it is in deep sleep, at which time its STM contains a NoOp chunk, i.e., a chunk with a NoOp gist and a high enough |weight| to keep all other chunks at bay. (See our answers to KM4 for discussions of Sleep and Dream processors.)

rCTM can have different *degrees of consciousnesses*. For one, its many processors are instrumental in developing Brainish-labeled sketches in rCTM’s world models. Involvement by processors like those for Smell, Vision, Hearing, Touch, raises the degree of conscious awareness. Even in deep sleep, however, a rCTM can still carry out tasks (utilizing unconscious communication between processors via links) but without attention and therefore without

## AI Consciousness is Inevitable

awareness. In addition, the degree of consciousness provided by the Brainish-labeled sketches is proportional – for as long as all else remains unchanged - to the |weight| of the chunk being attended to.

Faulty processors or faults in the competition tree can affect what gets into STM, hence affect both conscious attention and conscious awareness. For example, a faulty rCTM can exhibit blindsight, meaning it can do things that are normally done with conscious sight, but without having the *feeling* that it is sighted (Blum & Blum, 2022). This can happen if the Vision processor fails to get its chunks into STM (e.g., relevant branches in the Up-Tree are broken or the Vision processor fails to give high enough |weight| to its chunks) while previously formed links enable visual information to enter and trigger the Motor processor.<sup>46</sup>

*Different degrees of consciousness* already occur in a developing rCTM. As we have noted, an infant rCTM has only a very foggy world model which does not even include a sketch of itself. Sketches with annotated labels develop and become refined *gradually*. They are what rCTM is consciously aware of.

**KM4. For things that have the capacity for consciousness, what distinguishes the *state* of consciousness from being unconscious? Is there a simple on/off switch? How is this related to arousal, attention, awareness of one’s surroundings (or general responsiveness)?**

In the winner-take-all competition, chunks reach STM with probability proportional to (a *monotonic function* of) the chunk’s weight. If all chunks have zero weight, then chunks flit in and out of STM at random and so fast that rCTM loses anything remotely resembling sustained attention (like Robbie the Robot in *Forbidden Planet*). We view this as a *state of unconsciousness*. rCTM can get out of this state only when some processor creates a nonzero-weighted chunk. (See Appendix 7.2 for discussion of the competition to enter STM.)

Another unconscious state occurs when a *Sleep processor* generates a NoOp chunk (a chunk having a NoOp gist) that has a “sufficiently high” |weight|. A sufficiently high |weight| is one well

---

<sup>46</sup> In the human visual cortex, the dorsal stream of vision is unconscious; the ventral stream is conscious. Studies on blindsight suggest that communication via the dorsal stream may account for blindsight in visually impaired people (Tamietto & Morrone, 2016).



## AI Consciousness is Inevitable

above the weight of any other chunk. That prevents other chunks - including those from processors that interact with the outer world<sup>47</sup> - from having much chance to enter STM. (But note, even a high |weight| chunk can always be replaced by a chunk of even higher |weight|.)

When the |weight| of a Sleep processor's chunk drops a bit - but not enough to let input-output chunks enter STM - a rCTM's *Dream processor* can take over, enabling chunks of a dream to emerge. If the |weight| of the Sleep processor's chunks drop even further, rCTM wakes up.

The above are some of the ways rCTM can go from consciousness to unconsciousness and back. The Sleep processor's weight assigning mechanism is a kind of *on/off switch* for consciousness in rCTM.

In an unconscious state, rCTM is not aware of its surroundings, though it might be *aroused* by pangs of intense hunger, other pains, a very loud explosion, and so on. This occurs, for example, when these pangs, pains, and sounds overwhelm the Sleep processor with even larger |weight|.

### **KM5\*. What determines what we are conscious of at any moment?**

In rCTM, at every clock tick, *t*, there is exactly one chunk in STM. When a chunk is broadcast, rCTM pays *conscious attention* to that one chunk only. Chunks are purposely small (in a well-defined way). This ensures that all processors can focus on the same thought.

### **KM6\*. Why do some neural or cognitive operations go on consciously and others subconsciously? Why/how are some kinds of information permitted access to our conscious awareness while most are excluded?**

Since we have not defined "subconscious" in the rCTM, we will substitute "unconscious" (meaning "not conscious") for "subconscious" in this answer.

Operations within each LTM processor are done unconsciously, i.e., they do not go through STM. Communication between LTM processors via links is also unconscious. Such communication is much quicker than conscious communication that goes through STM.

Alison Gopnik's contention that "babies are more conscious than we are" (Gopnik, 2007) can be understood in terms of what we call *conscious attention*. In the *infant* rCTM, until links are

---

<sup>47</sup> Something like this can happen in total depression and catatonia in rCTM, and slow-wave (non-REM) sleep in humans.

## AI Consciousness is Inevitable

formed, all communication between processors is conscious, i.e., goes through STM. Then as processors form links, communication can go quickly through links, meaning unconsciously. This is what happens after the young rCTM learns to ride a bike.<sup>48</sup>

On the other hand, infant rCTM's MotW is considerably less developed than the adult rCTM's MotW. Hence phenomenal consciousness (what we call *conscious awareness*) is considerably less developed in the infant rCTM than in the adult.

**KM7\*. What distinguishes things that we are currently consciously aware of, from things that we *could be* consciously aware of if we turned our attention to them, from things that we could not be consciously aware of (that nevertheless play crucial roles in our cognition)?**

For rCTM to be consciously aware of a thing, call that thing **abc**, a chunk referring to **abc** must get into STM. Once it does, rCTM pays conscious attention to **abc**. But even conscious attention to **abc** does not make for conscious awareness of **abc**. For that, the chunk must reference a sketch in the MotW that is called or labeled with **abc**.

What things, though important for cognition, *cannot* enter consciousness? Here are a couple of answers from rCTM:

1. Things that must be done so quickly that the communication necessary to do the thing cannot go through STM. For example, rCTM must quickly swerve away from an oncoming car while riding its bike.
2. Things like **abc** whose doing would take away from a much more important thing like **xyz**. In that case, time permits only one of **abc** and **xyz** to be consciously attended to. If there is barely enough time to do one (and only one) of them, then rCTM is much more likely to be {unconscious of **abc** and conscious of **xyz**}.

---

<sup>48</sup> Humans learn to play ping pong consciously. In a ping pong tournament however, a player must let the unconscious take over, must insist that the conscious get out of the way. In swimming, repetition gives one's unconscious an opportunity to improve one's stroke, but it doesn't enable a new stroke to be acquired. That requires conscious attention. For example, the dolphin kick is weird and unnatural, but since it works for dolphins, it makes sense to simulate it, and that is done consciously at first. The unconscious then optimizes the constants.

## AI Consciousness is Inevitable

**KM8. Which systems are required to support conscious perception? Where is the relevant information represented? Is it all pushed into a common space or does a central system just point to more distributed representations where the details are held?**

In rCTM, conscious awareness is impossible without the MotWp (among others). Conscious attention *is* possible without the MotWp, but impossible without the broadcast station.

To support conscious perception, the relevant information is held in the MotW. For example, color is in the Color processor, smell is in the Smell processor, and so on. So, *in that sense*, information is *distributed* (even though we have stipulated that the sensory processors are part of MotWp.) When rCTM first sees and smells a rose, the MotWp creates and attaches the labels RED and SMELLS\_SWEET to the sketch of the rose in the MotW. At some point in time, RED and SMELLS\_SWEET become fused as a Brainish word or gist, and *in that sense*, information is *unified*. This can also occur in a processor that takes on the task of linking to the processors for color and smell.

**KM9\*. Why does consciousness feel unitary? How are our various informational streams bound together? Why do things feel like *\*our\** experiences or *\*our\** thoughts?**

At each clock tick, all LTM processors *simultaneously receive a global broadcast* of the conscious content (current chunk) in STM. That gives rCTM its sense of a *unitary* experience. Additionally, when the broadcasted chunk contains a gist that refers to a sketch with a *fused multimodal Brainish label*, that broadcasted information feels unitary. Information passed through *links* further *bind information* together.

If rCTM's conscious content refers to a thought or experience that MotWp has labeled SELF, rCTM will be consciously aware of that thought as its own. If that thought is also labeled FEELS, rCTM will not only *know* that the thought is its own, it will also *feel* that it is its own.

**KM10\*. Where does our sense of selfhood come from? How is our conscious self related to other aspects of selfhood? How is this *sense of self* related to actually *being a self*?**

Here again, world models, with their *learned* Brainish labels, determine rCTM's sense of self. The MotW's sketches are labeled with a variety of gists. For this question, the labels FEELS, SELF, and CONSCIOUS are particularly important. If all three labels are attached to a sketch of rCTM in the MotW and this sketch with labels is broadcast, then rCTM **FEELS** that its **SELF** is **CONSCIOUS**.

## AI Consciousness is Inevitable

Known pathologies occur when any one or more of these labels is missing, or when sketches are mislabeled.<sup>49</sup>

### **KM11. Why do some kinds of neural activity feel like something? Why do different kinds of signals feel different from each other? Why do they feel specifically like what they feel like?**

In rCTM, inputs from different sensors go to *different sensory processors*. Those different senses become incorporated in the MotW with *different Brainish labels*.

In the MotW, sketches of a red rose and a red fire engine will both get labeled RED. But these sketches will get many other labels as well. For example, the fire truck sketch likely gets the Brainish labels FIRE\_TRUCK and LOUD\_SIREN while the rose sketch does not. The rose sketch gets labeled “FEELS\_SILKY and SMELLS\_SWEET” while the fire truck does not.

The two referents are distinguished in the MotW, and with increasingly more Brainish labels, “feel specifically like what they feel like.” This distinction is the result of Brainish’s and MotW’s evolutionary development during rCTM’s lifetime.

Similar arguments apply to the distinct pleasures rCTM feels when seeing a red rose and when satisfied after when replenishing its fuel tank. See also our answer to the next question.

### **KM12\*. How do we become conscious of our own internal states? How much of our subjective experience arises from homeostatic control signals that necessarily have valence? If such signals entail feelings, how do we know what those feelings are about?**

In the Overview (section 2.3), we indicated how the infant rCTM would know it is hungry when a high |weight| negatively valenced chunk from its Fuel Gauge processor reaches STM and is broadcast from it.

The LOW\_FUEL chunk will trigger an actuator to connect rCTM’s fuel intake to a fuel source (in humans, the breast). Assuming it works, that will eventually result in a sketch of the fuel source (in the MotW) being labeled FUEL\_SOURCE and PLEASURE\_SOURCE. At the same time, the labels

---

<sup>49</sup> Some human examples of pathologies due to mislabeling include body integrity dysphoria (when SELF is missing from some body part), phantom limb syndrome (when an amputated arm is still labeled SELF), Cotard’s syndrome (when SELF and FEEL are missing from the representation of oneself in the MotW), anosognosia (when ERROR DETECTION is missing from the representation of oneself), paranoia (when a friend is labeled SPY), ....

## AI Consciousness is Inevitable

FEELS\_HUNGRY, INTAKES\_FUEL and FEELS\_PLEASURE will be attached to the sketch of rCTM while it is hungry and being fueled. A positive weight broadcast indicates that “rCTM feels pleasure while it gets fuel if it’s hungry.” This process is an example of homeostasis in rCTM and how rCTM becomes conscious of its own internal state.

The about-ness of those feelings come from the Brainish labels and sketches that evolve during rCTM’s lifetime. Thus FEELS\_PLEASURE could *be about* being connected to rCTM’s fuel source. FEELS\_HUNGER and FEELS\_PAIN will *be about* rCTM’s Fuel Gauge registering low fuel and the sketch of the CTM in the MotW being labeled LOW\_FUEL. See also our answer to the next question.

### **KM13. How does the about-ness of conscious states (or subconscious states) arise?**

**How does the system know what such states refer to? (When the states are all the system has access to).**

Conscious states in rCTM are broadcasted states. The MotW is *all* that rCTM knows about its (inner and outer) worlds, this includes its conscious states. Labeling the sketch of the state in the MotW and broadcasting this labeled sketch determines the about-ness of these states and how rCTM knows what such states refer to.

For example, suppose the conscious state is of the world in front of the CTM, which includes a staircase, and rCTM wants to know the effect of starting down the stairs. The MotWp predicts the effect of starting down the stairs through a simulation of that action in the MotW.<sup>50</sup> Say the CTM decides to go ahead. Responses from the world provide a Brainish label to the sketch (succinct description) that includes the prediction, the action, and the actual response to that action is the about-ness of the state.

When a chunk with gist referring to this state’s sketch and label is broadcast, all processors learn the about-ness of the current conscious state.

### **KM14. What is the point of conscious subjective experience? Or of a high level common space for conscious deliberation? Or of reflective capacities for metacognition? What adaptive value do these capacities have?**

A conscious subjective feeling is experienced when broadcasted chunks incorporate rich Brainish-labeled sketches in the MotW. These labels evoke the feelings that conscious awareness is about.

---

<sup>50</sup> This is similar to the kind of simulation that the MotW does in a dream sequence.

## AI Consciousness is Inevitable

Without these feelings, rCTM would not be compelled to act appropriately. A person who has pain and *knows* everything about it but lacks the ability to *feel* its agony has *pain asymbolia*. Such a person is not motivated to respond normally to pain, or to take care of themselves. Same for rCTM.

With conscious subjective feelings, reflective capacities enable rCTM to treat itself with all the planning tools it uses to treat other referents.

Global broadcasting, by focusing all processors on the same thing, creates a high level space for deliberation. It also enables conscious subjective experiences to be adaptive. By receiving each and every broadcast, all processors can contribute to the understanding of the broadcast and/or its solution.

For example, suppose the broadcast is a high |weight| negatively valenced processor scream, indicating a problem: “Hungry! Must fill the fuel tank.” To deal with the situation, the Navigation processor might contribute a choice of routes to local fuel stations. The Computation processor might compute how much fuel is required for each choice. The Weather processor might weigh in that one of the routes is blocked. Employing these suggestions to solve the problem can alter rCTM’s subjective experience of hunger to one of satisfaction and well-being.

**KM15\*. How does mentality arise at all? When do information processing and computation or just the flow of states through a dynamical system become elements of cognition and why are only some elements of cognition part of conscious experience?**

The question asks: “How does the capacity for intelligent thought come about?” rCTM is ideal for answering this question since at birth, all  $10^7$  processors are independent. In the case of the O2 pain, the first processors to come online – meaning they have sufficient weight to get their chunks to STM - are those having homeostatic importance like the Nociceptor Gauges (monitor pain), Fuel Gauge (monitors hunger), and so on, or have immediate access to the senses like vision, hearing, and so on. These processors help the MotWp to make and improve its predictions and world models. The next processors to come online are those that affect the activators, one of which cries for help. Then come processors that detect coincidences like: “This visual input and that auditory input seem to coincide.” This is the beginning of intelligent thought.

The rCTM model, unlike Baars’ GW, has no Central Executive. The *competition for conscious attention*, which replaces the Central Executive, gives rCTM much of its cognitive power. That competition efficiently considers *all* information submitted for consideration by its  $10^7$  processors. It allots to each idea a winning probability or share of consciousness (broadcast time)

## AI Consciousness is Inevitable

proportional to its estimated importance).<sup>51</sup> It enables processors to solve a problem even though rCTM does not know which processors have the interest, expertise or time to consider the problem. No Central Executive could have the knowledge or resources to do that (unless, of course, the Central Executive was itself an efficient competition process).

Some elements of cognition can be performed within a single LTM processor and transmitted through links. This processor doesn't need to search through an enormous data base for its information: it already has the information or knows where the necessary information is held. This is unconscious cognition. Processors that do need to search for the information may need to broadcast their need.<sup>52</sup> That broadcast begins the process of using consciousness to do cognition.

**KM16\*. How does conscious activity influence behavior? Does a capacity for conscious cognitive control equal “free will”? How is mental causation even supposed to work? How can the meaning of mental states constrain the activities of neural circuits?**

In rCTM, *conscious activity is intertwined with behavior*.

In rCTM, all LTM processors receive the broadcasted conscious content. Different processors have differing amounts of time to deal with that content. Of those that have time, some have a more reasonable idea how to deal with the broadcast than others. A broadcasted message that the Fuel Gauge is low can prompt one processor to try to conserve fuel, another to trigger a search for a source of fuel, and so on. A broadcast of danger may prompt rCTM to choose between fight, flight or freeze, each championed by a different processor.

Additionally, rCTM's *disposition* plays an important factor in the competition that selects which chunk will be globally broadcast and hence its behavior. (See Appendix 7.2 for more information on rCTM's competition and the influence of its disposition.)

As for “free will”, rCTM's ability to assess a situation, consider various possibilities, predict the consequences of each, and based on that make a decision (all *under resource constraints*) gives rCTM its feeling of “free will”. For example, imagine rCTM playing a game of chess. When and for as long as rCTM has to decide which of several possible moves to make, it knows it is “free” to choose

---

<sup>51</sup> This is something that tennis and chess tournaments do not provide. (See Appendix 7.2.)

<sup>52</sup> Like the processor that asks, “What her name?” (See footnote in section 2.1).

## AI Consciousness is Inevitable

whichever move it determines has the greatest utility for it. That is free will. See (Blum & Blum, 2022).

In rCTM, the MotW is fundamental to *mental causation*. To will an act in the world, the MotWp first simulates that action in the MotW, gauges whether or not to do it, if so gets it done, then looks to see if the act got accomplished in the world and if it predicted correctly.

For example, suppose the infant rCTM discovers that it can somehow move its left leg. It becomes aware through its sensors that “willing” the movement of that leg is successful. For comparison, it may discover that it cannot pick up a rock, Yoda style, with the power of thought. Moving the leg or lifting the rock can be “willed” by performing the action in the MotW. Sensors must verify if the act has been successful. If it has, that is mental causation.

As an example of how mental states can constrain activities, in our answer to KM4, we discussed how the Sleep processor generates a non-dreaming sleep state by raising its own |weight| so high that other chunks can’t reach STM. This shows how the sleep state *constrains* activity in rCTM’s Up-Tree.

-----

We’ve now come to the end of Kevin Mitchell’s questions. He ends his blog with the words, “**If we had a theory that could accommodate all those elements and provide some coherent *framework*<sup>53</sup> in which they could be related to each other – not for providing all the answers but just for asking sensible questions – well, that would be a theory of consciousness.**”

## 5 Summary and Conclusions

In this paper we have presented a brief overview of a simple formal machine model of consciousness, known here as rCTM. The Theoretical Computer Science perspective has influenced the design and definitions of rCTM, and conclusions we have drawn from the model.

Although rCTM is inspired by the simplicity of Turing’s formal model of computation and Baars’ global workspace (GW) architecture, our formalization is neither a Turing Machine nor a standard

---

<sup>53</sup> Italics ours.



## AI Consciousness is Inevitable

GW model. Its *consciousness* (access and phenomenal) depends on what is under the hood and it *having more than a global workspace*.

Importantly, rCTM also:

1. interacts with its outer world via input sensors and output actuators;
2. has the ability to construct models of its inner and outer worlds;
3. has a rich internal multimodal language, Brainish; and
4. constantly updates its states via predictive dynamics (cycles of prediction, testing, feedback and learning),

all while operating under resource limitations (time and space).

rCTM is not a model of the human or animal brain, nor is it intended to be. It is a simple formal *machine model* of consciousness. Nevertheless, at a high level, rCTM can exhibit phenomena associated with human consciousness (blindsight, inattention blindness, change blindness, body integrity identity disorder, phantom limb syndrome, ...), and aligns with and integrates those key features from main theories of consciousness that are considered essential for human and animal consciousness.

The rCTM model thus demonstrates the compatibility and/or complementarity of those theories., And because the rCTM is clearly buildable and arguably a basis for consciousness, it further supports (the credibility of) our claim that *a conscious AI is inevitable*.

Finally, the development of rCTM is a work in progress. While we have worked out many details of the model, there is much left to develop. More specifics will appear in our upcoming monograph.

Our goal is to explore the model as it stands, determine the good and the bad of it, and make no changes to it.

## 6 Acknowledgements

We are grateful to Michael Xuan for his immense encouragement, and to UniDT for their long-term support. We appreciate numerous discussions with and helpful critiques from friends and colleagues: Johannes Kleiner, Wanja Wiese, Ron Rivest, Hy Hartman, Alvaro Velasquez, Sergio Frumento. We thank Kevin Mitchell for posing sensible questions “that a theory of consciousness should be able to encompass” and his subsequent insightful comments and discussion with us.

### 7 Appendix

#### 7.1 A Brief History and Description of the TCS Approach to Computation

The theoretical computer science approach to computation starts with Alan Turing in the 1930's and focuses on the question, "What is computable (decidable) and what is not?" (Turing, 1937). Turing defined a simple formal model of computation, which we now call the Turing Machine (TM) and *defined* a function to be computable if and only if it can be realized as the input-output map of a TM. The formal definition of a TM (program) also provides a formal definition of the informal concept of algorithm.

Using his model, Turing proved properties (theorems) of computable functions, including the existence of universal computable functions (universal Turing machines) and the fact that some functions are not computable. The former foresees the realization of general purpose programmable computers; the latter that some problems cannot be decided even by the most powerful computers. For example, Turing shows there is no Turing machine (Turing computable function) that given the description of a TM  $M$  and an input  $x$ , outputs  $1$  if  $M$  on input  $x$  (eventually) halts, and  $0$  if not. This is known as the "halting problem" and is equivalent to Gödel's theorem on the undecidability of arithmetic.

But why should we believe the Church-Turing Thesis, suggested first in (Turing, 1937), that the TM embodies the informal notion of computability (decidability)? That's because each of a great many very different independently defined models of discrete computation, including TMs and Alonzo Church's *effective calculability* (Church, 1936), define exactly the same class of functions, the computable functions. In programming parlance, all sufficiently powerful practical programming languages are equivalent in that anyone can simulate (be compiled into) any other. The ensuing mathematical theory is often called the Theory of Computation (TOC).

In the 1960's, with the wider accessibility of computers, newly minted theoretical computer scientists such as Jack Edmonds (Edmonds, 1965) and Richard Karp (Karp, 1972), pointed out that resources matter. Certain problems that in principle were decidable, were seemingly intractable given feasible time and space resources. Even more, intractability seemed to be an intrinsic property of the problem, not the method of solution or the implementing machine. The ensuing sub-theory of TOC, which introduces resource constraints into what is or is not computable *efficiently*, is called Theoretical Computer Science (TCS).

## AI Consciousness is Inevitable

TCS focuses on the question, “What is or is not computable (decidable) given limited resources?” A key problem here is the deceptively simple “SAT problem”: Given a boolean formula  $\mathcal{F}$ , is it satisfiable, meaning is there a truth assignment to its variables that makes formula  $\mathcal{F}$  true? This problem is decidable. Here is a decision procedure: Given a boolean formula  $\mathcal{F}$  with  $n$  variables, systematically check to see if any of the  $2^n$  possible truth assignments makes the formula true. If yes, output **1**, otherwise output **0**. This brute force procedure takes exponential( $n$ ) time in general. But is the “SAT problem” tractable, meaning decidable efficiently, i.e., in polynomial( $n$ ) time? This is equivalent to the well-known  $P = NP?$  problem of (Cook, 1971), (Karp, 1972), (Levin, 1973).

The design of novel and efficient algorithms is a key focus of TCS.

Turning the table on its head, the ability to exploit the power of hard problems, problems that cannot be solved efficiently, has been a key insight of TCS. An example is the definition of pseudo-randomness (Hatami & Hoza, 2024). This ability to exploit the power of hardness is novel for mathematics.

### 7.2 The Probabilistic Competition for Conscious Attention and the Influence of Disposition on it

We tried to make the rCTM Up-Tree competition deterministic, but it turns out to *necessarily* be probabilistic: This is because any deterministic competition must be made increasingly complex to realize essential properties. For example, this is the case if we want chunks of near equal |weight| to have near equal chances of getting onto the STM stage.

Consider a deterministic competition which makes decisions based on a chunks weight. Suppose chunk A is pinned to weight 11, chunk B to weight 9, and all other chunks to weight 0. In the deterministic rCTM competition, A always wins and B never does. And as long as weights don’t change, rCTM remains totally unconscious of chunk B. In the probabilistic rCTM competition described below (with simple **f value** = intensity), chunk A wins with probability 11/20 while B win with probability 9/20. So, A and B each have roughly equal probability of winning a competition, and a new independent competition is begun at every clock tick. In that case, rCTM will likely become conscious of both.

In the probabilistic rCTM competition (a variant of the standard tennis or chess tournament) it is proved that a chunk wins the competition with probability proportional to (a *function* of) its weight. As a consequence, the winning chunk is independent of where processors are located! This is a property of rCTM’s probabilistic competition. (It is a property that would be difficult if not impossible to achieve in tennis and chess tournaments.)

## AI Consciousness is Inevitable

We now describe the probabilistic competition. First recall that a *chunk* is a tuple,

**<address, time, gist, weight, auxiliary information>**,

consisting of the address of the originating processor, the time the *chunk* was put into the competition, a succinct Brainish *gist* of information, a valenced weight (to indicate the importance/ value/ confidence the originating processor assigns its gist), and some auxiliary information.

For the probabilistic rCTM, the **auxiliary information** is a pair of numbers which we call **(intensity, mood)**.

*At the start of the competition*, each LTM processor enters a chunk into its leaf node with

**intensity = |weight| and mood = weight.**

In the probabilistic competition, each non-leaf node of the Up-Tree contains a *coin-toss neuron*.

The coin-toss neuron probabilistically chooses the *local winner* of the two (competing) chunks in the node's two children based on their **f values**. Here **f** is a function mapping chunks to non-negative real numbers. A *simple*, but natural **f** value maps chunks to their intensities.

If  $C_1$  and  $C_2$  are the two competing chunks, then the coin-toss neuron will choose  $C_i$  to be the local winner with probability  $f(C_i)/f(C_1)+f(C_2)$ .

The local winner will move up one level in a single clock tick. The first four parameters of this new chunk are the same as the local winner's. But *its intensity* is the sum of two competing chunks' intensities. Similarly for its *mood*. We call this the WINNER TAKE ALL policy!

Thus, as a chunk moves up the tree, the intensity never decreases. This is not (necessarily) the case for mood.

In this way, the winning chunk's auxiliary information at the end of the competition will contain the *sum of all submitted chunks' intensities* (|weights|) and *the sum of all submitted moods* (weights).

Hence, although at the start of the competition each processor has little idea about the other  $N \gtrsim 10^7$  chunks that are being submitted, the reception of the broadcasted winner provides useful information. In addition to providing the winning gist and its processor's address, the broadcasted

## AI Consciousness is Inevitable

information enables each LTM processor to quickly compare how *its submitted chunk's weight* (and hence *intensity* and *mood*) compares with the winner's.<sup>54</sup>

So, for example, in the competition with the simple  $f$  values mentioned above, each processor can easily determine the average of all submitted chunks' intensities and moods, by dividing the broadcasted intensity and mood by  $N$ .

More generally, consider an rCTM whose competition employs the following  $f$  value:

$$f(\text{chunk}) = \text{intensity} + d \cdot (\text{mood}) \quad \text{and} \quad -1 \leq d \leq +1.$$

Here  $d$  is a fixed constant called rCTM's *disposition*.

rCTM's *disposition* plays an important factor in the competition that selects which chunk will be globally broadcast, and hence rCTM's behavior.

If the disposition is  $d = 0$  we say rCTM is "level headed". In this case, the probability of an entered chunk "winning" the competition will depend on its  $|\text{weight}|$ , independent of valence.

If the disposition is  $d > 0$ , rCTM will be "upbeat" in the sense that positively valenced chunks will have a higher probability of winning than negatively weighted chunks of the same  $|\text{weight}|$ . If its disposition is  $d = +1$ , rCTM is manic, only positively valenced chunks can win the competition. The rCTM knows only what is positive in its life, as long as anything is positive.

If the disposition is  $d < 0$ , rCTM will be "downbeat". In the extreme, if  $d = -1$ , rCTM will be "hopelessly depressed". Only negatively valenced chunks can win the competition. There is no way out of this horrible state except with a reboot, i.e., to "shock" the system to get a less extreme disposition.<sup>55</sup>

---

<sup>54</sup> Question: Why does rCTM need a competition? Why not have each processor compute the probability of its chunk winning and then choose the winning chunk based on these probabilities? Answer: Again, "at the start of the competition each processor has little idea about the other  $N \geq 10^7$  chunks that are being submitted". Without this knowledge, the simple competition is an efficient way to ensure that each chunk gets into STM according to its probability based on  $|\text{weight}|$ s.

<sup>55</sup> In humans, electroconvulsive therapy (ECT) is used primarily for extreme depression. If  $d = +1$ , rCTM is in the manic state, and rCTM theory suggests that there too, a reboot is warranted. However, despite the fact that "ECT is a rapid and highly effective treatment of manic episodes", current guidelines only endorse ECT "for pharmacotherapy-resistant mania but often as second- or third-line treatment." (Elias, Thomas, & Sackeim, 2021).

## AI Consciousness is Inevitable

### References

- Baars, B. J. (1997). *In the Theater of Consciousness*. New York: Oxford University Press.
- Baars, Bernard J. (1997). In the Theater of Consciousness: A rigorous scientific theory of consciousness. *Journal of Consciousness Studies*, 4(4), 292-309.
- Baddeley, A. D., & Hitch, G. J. (1974). Working memory. In G. A. Bower (Ed.), *The Psychology of Learning and Motivation* (pp. 47-89). New York: Academic Press.
- Bancu, P., Bisnath, Burgess, Eakins, Gonzales, Saltzman, . . . Baptista. (2024). Revitalizing Attitudes Toward Creole Languages I. In A. H. Hudley, C. Mallinson, & M. Bucholtz, *Decolonizing Linguistic*. Oxford University Press.
- Berridge, K. C., & Kringelbach, M. L. (2015). Pleasure systems in the brain. *Neuron*, 86(3), 646-664. doi: 10.1016/j.neuron.2015.02.018.
- Block, N. (1995). On a confusion about a function of consciousness . *Brain and Behavioral Sciences*, 18(2), 227-247.
- Blum, A. (1995, July). Empirical support for winnow and weighted-majority algorithms: Results on a calendar scheduling domain. *Proceedings of the Twelfth International Conference on Machine Learning*, (pp. 64-72. <https://doi.org/10.1016/B978-1-55860-377-6.50017-7>).
- Blum, A., Hopcroft, J., & Kannan, R. (2015). *Foundations of Data Science*. Ithaca. Retrieved from <https://www.cs.cornell.edu/jeh/book.pdf>
- Blum, L., & Blum, M. (2022, May 24). A theory of consciousness from a theoretical computer science perspective: Insights from the Conscious Turing Machine. *PNAS*, 119(21), <https://doi.org/10.1073/pnas.21159341>.
- Blum, L., & Blum, M. (2023). A Theoretical Computer Science Perspective on Consciousness and Artificial General Intelligence. *Engineering*, 25(5), 12-16. <https://doi.org/10.1016/j.eng.2023.03.010>.
- Blum, M., & Blum, L. (2021, March). A Theoretical Computer Science Perspective on Consciousness. *JAIIC*, 8(1), 1-42. <https://doi.org/10.1142/S2705078521500028>.
- Blum, M., & Blum, L. (2022, December 24). *A Theoretical Computer Science Perspective on Free Will*. Retrieved from ArXiv: <https://doi.org/10.48550/arXiv.2206.13942>

## AI Consciousness is Inevitable

- Carney, J. (2020). Thinking avant la lettre: A Review of 4E Cognition. *Evolutionary Studies in Imaginative Culture*, 4(1), 77-90. <https://doi.org/10.26613/esic.4.1.172>.
- Chalmers, D. J. (1995). Facing Up to the Problem of Consciousness. *Journal of Consciousness Studies*, 2(3), 200-219.
- Church, A. (1936). A note on the Entscheidungsproblem. *J. of Symbolic Logic*, 1 (1936), 40-41, 1, 40-41.
- Clark, A. (2008, Jan). Pressing the Flesh: A Tension in the Study of the Embodied, Embedded Mind? *. JPhilosophy and Phenomenological Research*, 76(1), 37-59. <https://www.jstor.org/stable/40041151>.
- Clark, A. (2015). Embodied prediction. In T. Metzinger, & J. Windt, *Open Mind*. Frankfurt am Main: MIND Group.
- Clark, A., & Chalmers, D. (1998, January). The Extended Mind. *Analysis*, 58(a), 7-19.
- Cleeremans, A. (2014). Prediction as a computational correlate of consciousness. *International Journal of Anticipatory Computing Systems*, 29, 3-13.
- Cook, S. A. (1971). The complexity of theorem-proving procedures. *Proceedings of the third annual ACM symposium on Theory of computing*, (pp. 151-158. <https://doi.org/10.1145/800157.805047>).
- Damasio, A. (1994). *The Feeling of What Happens*. NY, NY: Harcourt, Brace and Co.,
- Dehaene, S. (2014). *Consciousness and the Brain: Deciphering How the Brain Codes Our Thoughts*. New York: Viking Press.
- Dehaene, S., & Changeux, J. P. (2005, April 12). Ongoing Spontaneous Activity Controls Access to Consciousness: A Neuronal Model for Inattentional Blindness. *PLoS Biol*, 3(5), <https://doi.org/10.1371/journal.pbio.0030141>.
- Dehaene, S., & Changeux, J. P. (2011, April 28). Experimental and theoretical approaches to conscious processing. *Neuron*, 70(2), 200-227. DOI: 10.1016/j.neuron.2011.03.018.
- Dehaene, S., & Naccache, L. (2001, April). Towards a cognitive neuroscience of consciousness: basic evidence and a workspace framework. *Cognition*, 79(1-2), 1-37. [https://doi.org/10.1016/S0010-0277\(00\)00123-2](https://doi.org/10.1016/S0010-0277(00)00123-2).

## AI Consciousness is Inevitable

- Dehaene, S., Lau, H., & Kouider, S. (2017, Oct 27). What is consciousness, and could machines have it? *Science*, *58*(6362), 486-492. doi: 10.1126/science.aan8871.
- Edelman, G. M. (2006, Summer). The Embodiment of Mind. *Daedalus*, *135*(3), 23-32.  
<https://www.jstor.org/stable/20028049>.
- Edmonds, J. (1965). Paths, trees, and flowers. *Can. J. Math.*, *17*, 449–467. doi:10.4153/CJM-1965-045-4.
- Elias, A., Thomas, N., & Sackeim, H. A. (2021). Electroconvulsive Therapy in Mania: A Review of 80 Years of Clinical Experience. *American Journal of Psychiatry*, *178*(3), 229-239. doi: 10.1176/appi.ajp.2020.20030238.
- Farisco, M., Evers, K., & Changeux, J.-P. (2024, April 18). *Is artificial consciousness achievable? Lessons from the human brain*. Retrieved June 2024, from arXiv:  
<https://arxiv.org/abs/2405.04540>
- Friston, K. (2005, April 29). A theory of cortical responses. *Phil. Trans. R. Soc. B*, *360*, 815-836. doi:10.1098/rstb.2005.1622.
- Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature reviews neuroscience*, *11*(2), 127-138. <https://doi.org/10.1038/nrn2787>.
- Ginsburg, S., & Jablonka, E. (2019). *The Evolution of the Sensitive Soul*. Cambridge, Massachusetts, US: MIT Press.
- Gopnik, A. (2007, December). Why babies are more conscious than we are. *Behavioral and Brain Sciences*, *30*(5-6), 503-504. <https://doi.org/10.1017/S0140525X0700283X>.
- Grahek, N. (2001; 2007). *Feeling Pain and Being im Pain*. Universitat Oldenburg; MIT Press (2nd edition).
- Graziano, M. S., Guterstam, A., Bio, B., & Wilterson, A. (2020, May-June). Toward a standard model of consciousness: Reconciling the attention schema, global workspace, higher-order thought, and illusionist theories. *Cognitive Neuropsychology*, *37*(3-4)(3-4), 155-172. Retrieved from doi:10.1080/02643294.2019.1670630
- Hatami, P., & Hoza, W. (2024). Paradigms for Unconditional Pseudorandom Generators . *Foundations and Trends® in Theoretical Computer Science*, *16*(1-2), 1-210. <http://dx.doi.org/10.1561/0400000109>.



## AI Consciousness is Inevitable

- Hohwy, J., & Seth, A. (2020). Predictive processing as a systematic basis for identifying the neural correlates of consciousness (preprint). *PsyArXiv*, psyarxiv.com/nd82g.
- Humphrey, N. (2023). *Sentience: The Invention of Consciousness*. Cambridge, Massachusetts, US: MIT Press.
- Karp, R. M. (1972). *Reducibility Among Combinatorial Problems*. (R. E. Miller, & J. W. Thatcher, Eds.) New York: Plenum.
- Lee, T. S., & Mumford, D. (2003). Hierarchical Bayesian inference in the visual cortex. *Journal of the Optical Society of America, Optics, image science and vision*, 20(7), 1434-1448.
- Lenharo, M. (2024, January 18). Consciousness - The future of an embattled field (The consciousness wars: can scientists ever agree on how the mind works? ). *Nature*, 625, 438-440. doi:10.1038/d41586-024-00107-7.
- Levin, L. A. (1973). Universal Sequential Search Problems. *Probl. Peredachi Inf.*, 9(3), 115-116.
- Li, F.-F. (2023). *The Worlds I See*. New York: Flatiron Books: A Moment of Lift Book.
- Liang, P. P. (2022, April 14). *Brainish: Formalizing A Multimodal Language for Intelligence and Consciousness*. Retrieved from ArXiv.
- Libet, B. (1985). Unconscious cerebral initiative and the role of conscious will in voluntary action. *Behavioral and Brain Sciences*, 8(4), 529-539.
- López-Barroso, D., & de Diego-Balaguer, R. (2017). Language Learning Variability within the Dorsal and Ventral Streams as a Cue for Compensatory Mechanisms in Aphasia Recovery. *Front. Hum. Neurosci.*, 11(476), doi: 10.3389/fnhum.2017.00476.
- Mashour, G. A., Roelfsema, P., Changeux, J.-P., & Dehaene, S. C. (2020, March 4). Conscious Processing and the Global Neuronal Workspace Hypothesis. *Neuron*, 195(5), 776-798. doi: 10.1016/j.neuron.2020.01.026.
- Maturana, H., & Varela, F. (1972). *De Maquinas y Seres Vivos, Autopoiesis: La organizacion de lo vivo*. Santiago de Chile: Editorial Universitaria, S. A.
- Maturana, H., & Varela, F. (1980). *Autopoiesis and Cognition: The Realization of the Living*. Boston: Reidel.

## AI Consciousness is Inevitable

- McClelland, J. L., & Rumelhart, D. E. (1981). An interactive activation model of context effects in letter perception: I. An account of basic findings. *Psychological Review*, *88*(5), 375-407. <https://doi.org/10.1037/0033-295X.88.5.375>.
- McWhorter, J. H. (1998, December). Identifying the Creole Prototype: Vindicating a Typological Class. *Language*, 788-818. <https://www.jstor.org/stable/417>.
- McWhorter, J. H. (2008). *Defining Creole*. USA: Oxford University Press.
- Miller, G. A. (1956). The Magical Number Seven, Plus or Minus Two: Some Limits on our Capacity for Processing Information. *Psychological Review*, *63*(2), 81-97. <https://doi.org/10.1037/h0043158>.
- Miller, M., Clark, A., & Schlicht, T. (2022). Editorial: Predictive Processing and Consciousness. *Rev.Phil.Psych.*, *13*, 797-808. <https://doi.org/10.1007/s13164-022-00666-6>.
- Mitchell, K. (2023, September 18). *What questions should a real theory of consciousness encompass?*. Retrieved February 2024, from Wiring the Brain: <http://www.wiringthebrain.com/2023/09/what-questions-should-real-theory-of.html>
- Nagel, T. (1974, October). What Is It Like To Be a Bat? *Philosophical Review*, *83*(4), 435-450. <https://doi.org/10.2307/2183914>.
- Parr, T., Da Costa, L., & Friston, K. 2. (2019). Markov blankets, information geometry and stochastic thermodynamics. *Phil.Trans.R. Soc.*, <http://dx.doi.org/10.1098/rsta.2019.0159>.
- Perani, D., Saccuman, M. C., Scifo, P., Awander, A., Spada, D. B., & (, e. a. (2011). Neural language networks at birth. *Proc. Natl. Acad. Sci.*, *108*, 16056-16061. doi: 10.1073/pnas.1102991108.
- Philip Gerrans. (2024). Pain suffering and the self. An active allostatic inference explanation. *Neuroscience of Consciousness*, *2024*(1), <https://doi.org/10.1093/nc/niae002>.
- Rowlands, M. J. (2010). *The New Science of the Mind From Extended Mind to Embodied Phenomenology*. Cambridge, MA, US: The MIT Press.
- Sergent, C., & Dehaene, S. (2005, July-November). Neural processes underlying conscious perception: experimental findings and a global neuronal workspace framework. *J Physiol Paris*, *98*(4-6), 374-384. doi: 10.1016/j.jphysparis.2005.09.006.

## AI Consciousness is Inevitable

- Seth, A. K. (2015). The Cybernetic Bayesian Brain - From Interoceptive Inference to Sensorimotor Contingencies. In T. Metzinger, & J. M. Windt, *Open MIND* (p. doi: 10.15502/9783958570108 23 | 24). Frankfurt am Main: MIND Group.
- Shanahan, M. (2005). Global Access, Embodiment, and the Conscious Subject. *Journal of Consciousness Studies*, 12(12), 46-66.
- Sigal, U.-K. (2022, May 4). Editorial: Simple and Simplified Languages. *Frontiers in Communication*, 7, <https://doi.org/10.3389/fcomm.2022.910680>.
- Solms, M. (2019). The Hard Problem of Consciousness and the Free Energy Principle. *Front. Psychol.*, 9(2714), doi: 10.3389/fpsyg.2018.02714.
- Solms, M. (2021). *The Hidden Spring: A Journey to the Source of Consciousness*. New York, NY, US: W. W. Norton and Company.
- Solms, M., & Friston, K. (2018). How and Why Consciousness Arises: Some Considerations from Physics and Physiology. *Journal of Consciousness Studies* 25 (5-6):202-238. *Journal of Consciousness Studies*, 25(5-6), 202-238.
- Storm, J., Klink, P., Aru, J., Senn, W., Goebel, R., Pigorini, A., . . . Pennartz, C. (2024, May 15). An integrative, multiscale view on neural theories of consciousness. *Neuron*, 112(10), 1531-1552, doi: 10.1016/j.neuron.2024.02.004.
- Tamietto, M., & Morrone, M. (2016, Jan 25). Visual Plasticity: Blindsight Bridges Anatomy and Function in the Visual System. *Current Biology*, 26(2), 70-73. <https://doi.org/10.1016/j.cub.2015.11.026>.
- Thompson, E. (2007). *Mind in Life: Biology, Phenomenology, and the Sciences of Mind*. Harvard University Press.
- Tononi, G. (2004, November 2). An information integration theory of consciousness. *BMC Neuroscience*, 5(42), 42-72. doi: 10.1186/1471-2202-5-42.
- Tononi, G., & Koch, C. (2015, May 19). Consciousness: here, there and everywhere? *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 370 (1668), <https://doi.org/10.1098/rstb.2014.0167>.

## AI Consciousness is Inevitable

- Turing, A. M. (1937). On Computable Numbers, with an Application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, 2(42), 230-265.  
<https://doi.org/10.1112/plms/s2-42.1.230>.
- VanRullen, R., & Kanai, R. (2021, May 14). Deep learning and the Global Workspace Theory. *Trends in Neurosciences*, 44(9), 692-704. doi: 10.1016/j.tins.2021.04.005.
- Varela, F., Thompson, E., & Rosch, E. (1991). *The Embodied Mind*. Camb, MA: MIT Press.
- von Helmholtz, H. (1866; 1962). *Treatise on physiological optics* (Vol. 3). (J. Southall, Ed.) New York, NY: Dover Publication.
- Wiese, W. T.-2. (2020, July 11). The science of consciousness does not need another theory, it needs a minimal unifying model. *Neuroscience of Consciousness*, 2020(1),  
<https://doi.org/10.1093/nc/niaa013>.
- Wolfe, J. M. (1998). Visual memory: What do you know about what you saw? 8(9), 303-304.
- Zacks, O., & Jablonka, E. (2023, September 13). The evolutionary origins of the Global Neuronal Workspace in vertebrates. *Neuroscience of Consciousness*,  
<https://doi.org/10.1093/nc/niad020>.