

International Journal of Machine Consciousness

Vol. 6, No. 2 (2014) 1–14

© World Scientific Publishing Company

DOI: 10.1142/S1793843014001316



A Mechanistic Theory of Consciousness

Michael S. A. Graziano* and Taylor W. Webb†

*Department of Psychology, Princeton University,
Princeton, NJ 08544, USA*

**graziano@Princeton.edu*

†taylor.w.webb@gmail.com

Recently we proposed a theory of consciousness, the attention schema theory, based on findings in cognitive psychology and systems neuroscience. In that theory, consciousness is an internal model of attention or an “attention schema”. Consciousness relates to attention in the same way that the internal model of the body, the “body schema”, relates to the physical body. The body schema is used to model and help control the body. The attention schema is used to model and help regulate attention, a data-handling process in the brain in which some signals are enhanced at the expense of other signals. We proposed that attention and the attention schema co-evolved over the past half-billion years. Over that time span, the attention schema may have taken on additional functions such as promoting the integration of information across diverse domains and promoting social cognition. This paper summarizes some of the main points of the attention schema theory, suggests how a brain with an attention schema might conclude that it has a subjective awareness, and speculates that the same basic properties can be engineered into machines.

Keywords: Attention; internal model; awareness.

1. Introduction

Recently we proposed a theory of consciousness, the attention schema theory, based on findings in cognitive psychology and systems neuroscience. We argued that the theory helps to make sense of a large body of experimental work [Graziano, 2013, 2014; Graziano and Kastner, 2011; Kelly *et al.*, 2014]. If the theory is correct, it suggests that subjective experience is scientifically understandable, mechanistic, and can be artificially reconstructed. In the theory, consciousness is more than a philosophical flourish. It is one of the tools that brains use to process information. Neuroscientists will never fully understand how the brain works without understanding consciousness, and engineers will never build fully capable computers without designing them into some version of the same tool. Although the attention schema theory was formulated from the perspective of psychology and neuroscience, it might be of interest in other areas of expertise. The purpose of the present paper

2 M. S. A. Graziano & T. W. Webb

1 is to communicate some of the main points of the theory to an audience outside
2 neuroscience.

3 Consciousness is difficult to study because of its diverse connotations. To some
4 people, consciousness is the sum total of a person's memories. To others, it is an
5 awareness of oneself at any moment in time. Many researchers focus on the qualia
6 of sensory events such as color or touch. Some study altered states of consciousness
7 including dreaming or meditation. Others study pathologies of consciousness caused
8 by brain damage. All of these aspects of consciousness are legitimate topics of
9 study.

10 The approach to consciousness taken here, however, is focused on the central piece
11 of the puzzle. How does the brain become aware of anything at all, whether it is
12 memory, self, or sensory event? What is subjective experience? Not all information in
13 the brain reaches awareness. Most of it does not. What makes the difference between
14 merely *processing* information and being *aware* of it? The following sections outline
15 the attention schema theory, introducing it through an evolutionary perspective. The
16 paper provides only a cursory summary of some aspects the theory. A more complete
17 exposition is provided in the book *Consciousness and the Social Brain* [Graziano,
18 2013].

19 20 21 **2. The Evolution of Attention**

22 To explain the attention schema theory, it is necessary to begin with the process of
23 attention and the mechanisms by which it is controlled in the brain. Although these
24 mechanistic issues may seem far removed from the more ethereal issues of awareness
25 and subjective experience, the link between attention and awareness will hopefully
26 become clear in the following sections.

27 The word "attention" is used in many ways both colloquially and scientifically.
28 Here we use a specific, neuroscientific definition [Beck and Kastner, 2009; Desimone
29 and Duncan, 1995]. Attention is a selection process by which some signals in the brain
30 are enhanced in strength at the expense of other, competing signals. The boosted
31 signals have a bigger impact on downstream systems. Those signals are more deeply
32 processed, more likely to be stored in memory for later use, and more likely to alter
33 behavioral output. In this definition, attention is a data-handling method.

34 The earliest neural networks may have lacked anything like attention. For
35 example, hydras appear to have an undifferentiated nerve net incapable of selective
36 signal enhancement. Hydras may have branched from other animals about 600
37 million years ago (MYA), though that number is not certain [Budd, 2008]. Nervous
38 systems that use some form of selective signal enhancement can be found in almost all
39 other animals that have been studied including crabs, flies, birds, and people [e.g.,
40 Barlow and Fraioli, 1978; Beck and Kastner, 2009; Mysore and Knudsen, 2013; van
41 Swinderen, 2012]. These phyla and classes of animals are thought to have diverged
42 from each other in the late Cambrian during the so-called Cambrian explosion,
around 550–520 MYA. It is therefore a reasonable guess that the earliest forms of

AQ: Pls check
should it be
520-550.

1 attention evolved roughly between 600 and 520 MYA. Attention then presumably
2 increased in complexity and sophistication in the past half billion years of evolution.

3 Visual attention in humans and monkeys is the most heavily studied example of
4 attention [Beck and Kastner, 2009; Desimone and Duncan, 1995]. In the primate
5 visual system, attention is many-layered. Competition among signals occurs within
6 and between multiple layers of processing, including subcortical nuclei and many
7 interconnected cortical visual areas. The competition is also biased or influenced by
8 signals that impinge on the visual system. For example, if you are looking at a pile of
9 change on the table, the visual representation of a dime might rise in signal strength
10 and temporarily win the competition. One way the dime's visual representation
11 might be boosted is if light sparkles from the dime, providing what is termed a
12 bottom-up bias. A second way the dime's representation might be boosted is if you
13 are engaging a cognitive process to find dimes, providing what is termed a top-down
14 bias. This complicated interaction of competing signals and biasing signals results in
15 a constantly shifting attentional state in which one or another visual representation
16 wins the competition of the moment and is more fully processed.

17 18 **3. The Body Schema and the Attention Schema**

19 In the roughly half-billion-year timespan during which attention evolved, the brain
20 presumably evolved ever more sophisticated mechanisms to control attention. The
21 attention schema theory focuses on one aspect of this regulation of attention. To
22 control something, it is useful to have a model or simulation of the thing to be
23 controlled. The usefulness of an internal model is now a generally recognized principle
24 of control engineering [Franklin *et al.*, 1989; Jacobs, 1993].

25 For example, the general wants to control his army. To help, he has a model army
26 of plastic men and tanks on a map. The model is not very accurate but helps in
27 keeping track and making predictions. Indeed a crucial aspect of a control model is
28 that it does not need to be perfectly accurate. It can be a cartoonish, approximate
29 depiction and still provide benefit to the control system.

30 A good example of a control model constructed by the brain is the body schema, or
31 internal model of the body. It is worth outlining some of the key features of the body
32 schema in detail because of its close relationship to the sense of self and consciousness.

33 Regions of the brain that span the somatosensory system, the visual system, and
34 the motor system, integrate many sources of information to construct an internal
35 model or simulation of the body [Graziano and Botvinick, 2002; Hwang and Shad-
36 mehr, 2005; Kawato, 1999; Wolpert *et al.*, 1995]. That model is constantly updated.
37 It keeps track of body segments, their sizes, shapes, joint angles, speed, force, the
38 tension on muscles, and other properties. The model can also help to make predic-
39 tions a few seconds into the future.

40 The body schema is notoriously inaccurate in two ways.

41 First, the body schema lacks physical detail. It lacks information on the specific
42 bone structure inside the body, on muscle attachment points and wrapping geometry,

1 on how the proteins myosin and actin bind and pull against each other to produce
2 muscle force, and so on. The body schema contains no detailed physical or mechanistic
3 information. It is a surface model. It depicts the surface shape of the body and a
4 few need-to-know items such as the overall hinged structure of the limbs and the
5 movement of joints.

6 Imagine that an outer space alien discovers humans but lacks access to a body for
7 dissection. The alien, however, has a brain-reading device that can read the infor-
8 mation contained in the body schema. The alien scientist foolishly thinks he can use
9 the body schema to inform himself about the actual human body. Alas the alien
10 arrives at some peculiar conclusions. He concludes that the human body is magical. It
11 is magical in this sense — it can move in elaborate ways, but contains no internal
12 mechanism or structure to support that movement. That is how the body schema
13 describes the body. The body schema is intrinsically inaccurate, like a cartoon sketch.

14 But more than that, the body schema sometimes makes outright mistakes. A
15 person's arm can be in one position and the body schema can register it in a different
16 position. Dissociations between the body and the body schema are quite easy to
17 produce and form the basis of many standard somatosensory illusions [e.g., Botvinick
18 and Cohen, 1998; Lackner, 1988].

19 Why does the brain have such a sloppy model of the body? The answer is pre-
20 sumably a balance between cost and benefit. It takes processing time and energy, as
21 well as neuronal space in the brain, to compute a body schema. To optimize survival,
22 the brain needs something that can be computed fast and on the fly. It is adaptive to
23 have a quick and dirty model as long as it is good enough to get by most of the time.

24 Many of the same principles evident in the body schema are theoretically trans-
25 ferable to an attention schema. Because a brain has a need to control its own
26 attention, theoretically it ought to construct a model of attention, or an attention
27 schema. That model should be a constantly updated description of what attention is,
28 what it means for a brain to attend to something, what the possible consequences of
29 attention are, and what signals in particular are the focus of attention at the moment.
30 That model is likely to be quick and dirty, lacking any detail about the neuronal
31 mechanism of attention, and sometimes flat out wrong, but nonetheless useful as a
32 rough model of the brain's state of attention. In the next section, we explore the
33 psychological implications of an internal model of attention and how it may relate to
34 subjective awareness.

35 36 **4. Properties of the Attention Schema**

37
38 We suggest that the attention schema gradually co-evolved with attention over the
39 last half-billion years. Presumably the attention schema began as something quite
40 simple and then grew in sophistication. Perhaps simple forms of an attention schema
41 are present in flies or sea slugs. But to understand the attention schema from a
42 psychological perspective, it is useful to consider a type of animal with a more
complex brain that evolved more recently.

1 Suppose a monkey looks at and attends to a banana. Again, by “attention”, what
2 is meant here is something mechanistic. Visual signals related to the banana win a
3 competition in the brain and rise in strength. The stronger signals then drive
4 downstream processes such as memory, response choice, and the sensory guidance of
5 behavior.

6 If the attention schema theory is right, however, the monkey’s brain does more
7 than pay attention to the banana. It also constructs a schematic model of that state
8 of attention. The model would require the following three chunks of information.

9 First, the brain must construct a model of the banana including information on its
10 color, three-dimensional shape, location in space, and other object-defining proper-
11 ties. This model is probably mostly constructed in the visual system.

12 Second, the brain must construct a model of the monkey. Perhaps that self-model
13 is partly the body schema.

14 Third, the brain must construct a model of the specific relationship between
15 subject and object, a model of attention itself.

16 In this theory, the monkey’s brain constructs a large, multi-part, internal model
17 that says in effect, “There is a me, there is a banana in front of me, and in specific I am
18 paying attention to that banana.” The internal model of attention must link together
19 something like that information.

20 A monkey has some capacity for higher cognition. When his higher cognition
21 receives information from that internal model, what does it learn? Cognition is only
22 as well-informed as the internal models that feed into it. It can do no better than that.
23 In a sense, cognition is captive to the brain’s internal models. Higher cognition is like
24 the space-alien scientist noted in the last section, the one that gains information
25 about the physical body only by accessing the incomplete information in the body
26 schema, and therefore mistakenly concludes that the body is magical. The monkey’s
27 higher cognition gains information about the state of attention only by accessing the
28 incomplete information in the attention schema.

29 The attention schema would certainly not describe attention in a physically
30 accurate way. The model would lack any of the mechanistic details of neurons and
31 signal competition. The monkey has no need to know that it has neurons and signals,
32 synapses or neurotransmitters. Instead the model would contain sketchy and
33 superficial information about attention. It would describe attention as a magical
34 state of knowing. Here, we mean “magical” in the sense used in the previous section:
35 A process that lacks any physical or mechanistic basis. The model would depict a
36 state of knowing without any physical basis for that knowing.

37 The model would depict that magical state of knowing as hovering inside the
38 body. It is a part of the monkey’s own self, wedded to his body schema. The model
39 would also depict some of the basic implications of that magical state of knowing: It
40 implies an ability to choose to act on the banana, and an ability to remember the
41 banana for future reference.

42 An attention schema would depict a mental possession or subjective *experience* of
the banana. It is useful to keep in mind the meaning of the word “subjective”. There

1 is a subject, a me. There is an object, the banana. And there is a relationship between
2 the two: The subject has mental possession of the object and thus is enabled to act in
3 certain ways with respect to the object.

4 When that monkey's higher cognition introspects, or accesses the data in that
5 internal model, the data informs it that there is a self and the self has a subjective
6 awareness, or experience, of the banana in front of it. The monkey's cognition has no
7 means to doubt this information. Nothing tells it that this information comes from an
8 inner construct. Nothing tells it that the construct is a cartoonish depiction of
9 something else. Nothing tells it that it is being fed any information at all. Higher
10 cognition learns only that subjective experience exists, is here, is inside, and has
11 possessed that banana.

12 The monkey is aware of the banana.

13 The theory is of course not specific to bananas. It works as well for a sound or a
14 touch, a memory or a thought. The monkey attends to item X. The monkey also
15 constructs an internal model of that state of attention. If higher cognition accesses
16 that internal model, it is informed that there is a self and the self has a subjective
17 awareness of X.

18 This account of awareness arguably has a certain inevitability to it. Brains
19 engage in attention. To control attention, in control theory, there ought to be an
20 internal model of it, or an attention schema. That attention schema would neces-
21 sarily leave out the physical details. It would depict a state of *knowing* that is
22 non-physical, without mechanism. And higher cognition would be captive to that
23 internal model. The creature would be certain that it has subjective awareness and
24 would have no basis for understanding the true source of that certainty. The theory
25 explains how a brain can arrive at the conclusion that it is aware of something
26 without even knowing that it has arrived at a conclusion or that the conclusion
27 derives from computation. This account is in many ways similar to the account
28 of Gazzaniga [1970] in which awareness is a self-explanatory narrative. It is also
29 similar to the account of Dennett [1992] in which ineffable experience is replaced by
30 computation.

31 In the attention schema theory, awareness is not an illusion. It is better described
32 as a caricature. A caricature is a distorted depiction of something real. The process of
33 attention does physically exist. The brain's model of it, however, is not entirely
34 accurate, and therefore introspection gives us a distorted understanding of attention
35 that we report as an ethereal awareness.

37 **5. The Relationship Between Awareness and Attention**

38 If the theory is correct, then awareness and attention should relate to each other in
39 the following three ways.

40 First, awareness and attention should usually covary. If you are attending to
41 something, then in most circumstances you should also be aware of it. This match
42 between awareness and attention is indeed usually present [Posner, 1994; Merikle and

1 Joordens, 1997; Mack and Rock, 1998; Mole, 2008; De Brigard and Prinz, 2010;
2 Prinz, 2011].

3 Second, awareness should differ from attention in certain key ways. Just as the
4 body schema can sometimes become misaligned from the body due to inaccuracies
5 inherent in any internal model, awareness should sometimes become misaligned from
6 attention. It should be possible to pay attention to something by all physiological
7 measures and yet fail to be aware of it. Many studies have now confirmed that indeed
8 it is possible to pay attention to an item and yet have no reportable awareness of it
9 [Baars, 1997; McCormick, 1997; Kentridge *et al.*, 1999; Lambert *et al.*, 1999; Ivanoff
10 and Klein, 2003; Lamme, 2003; Woodman and Luck, 2003; Kentridge *et al.*, 2004;
11 Ansorge and Heumann, 2006; Jiang *et al.*, 2006; Koch and Tsuchiya, 2007; Mele
12 *et al.*, 2008; Mulckhuysen and Theeuwes, 2010; van Boxtel *et al.*, 2010]. It may seem
13 counter-intuitive to pay attention to something and yet be unaware of it. But
14 attention is a mechanistic process in the brain, like the regulation of blood flow or the
15 growth of new synapses. It is a process of signal enhancement. Awareness, in contrast,
16 is in the form of knowledge that is represented in the brain and can at least sometimes
17 be reported. Awareness acts, in effect, like the brain's sometimes-wrong knowledge of
18 its state of attention.

19 Third, when the brain attends to an item and yet is not aware of it, according to
20 the theory, the brain has a temporarily faulty internal model of its attentional state.
21 Therefore, the control of attention should suffer. By analogy, when the brain lacks a
22 clear internal model of the arm, the control of the arm is compromised. The arm may
23 be difficult to move to a new position or difficult to maintain in one position against
24 external perturbations [Scheidt *et al.*, 2005]. Just so, if you are attending to a visual
25 stimulus but unaware of it, your attention may be harder to disengage from the
26 stimulus, or may be unduly influenced by inconsequential features of the stimulus.
27 This third hypothesis about the relationship between awareness and attention —
28 that in the absence of awareness, the control of attention should act as though it has
29 lost its internal model — is one of the most crucial predictions of the theory. We are
30 currently testing it in human psychophysical studies.

31 32 33 **6. Integration of Information**

34 Many scholars believe that a defining feature of consciousness is its integration of
35 information across different domains [e.g., Baars, 1983; Crick and Koch, 1990;
36 Damasio, 1999; Engel and Singer, 2001; Newman and Baars, 1993; Schiff, 2008;
37 Tononi, 2008]. Although this integration of information is not the central contention
38 of the attention-schema theory, the theory is nonetheless compatible with the inte-
39 gration hypothesis. Indeed, the theory may provide a simple explanation for why
40 consciousness tends to be integrative.

41 The brain constructs models, or simulation, or updatable descriptions, of things in
42 the real world. Those models themselves are made of smaller components. For
example, for the visual system to construct a model of a red apple, it must link

1 together its model of the color red with its model of other stimulus features such as
2 the shape or movement of the apple. This is integration of information. Color,
3 however, is domain specific. It is not typically bound to information in other domains.
4 Unless you have a condition called synesthesia, you do not literally see sounds as
5 colored, see emotions as colored, or see mathematical thoughts as colored. Color can
6 be linked to other visual information, but not typically to information outside the
7 visual domain. It does not serve as a useful domain-general hub — a model to which
8 models of many other kinds can be linked.

9 But in the attention schema theory, the brain does construct a model that is
10 domain general. Attention is relevant to almost all domains of information processed
11 in the brain — to vision, sound, a touch on the skin, emotion, thought, memory, or
12 whatever the signals may be to which you are attending. In the attention schema
13 theory, the brain constructs a model of attention and links it to a model of the
14 attended item. That model of attention, the attention schema, is therefore an inte-
15egrative hub. It is domain-general — a model that is linkable to almost any category
16 of information.

17 Evolution is opportunistic. Sometimes a trait that evolves for one function takes
18 on other functions. Perhaps the attention schema evolved first as a way of helping to
19 control one's attention. We propose that a second obvious adaptive advantage of an
20 attention schema is its ability to link information across domains. In this theory,
21 awareness evolved initially as part of the control mechanism for attention and then
22 allowed for an increase in intelligence by promoting domain-general integration of
23 information.

24 25 **7. Social Cognition**

26 Over the half-billion years of its evolution, the attention schema may have taken on
27 many adaptive functions. We proposed that it was gradually modified to model,
28 monitor, and predict the attentional states of other animals [Graziano, 2013, 2014;
29 Graziano and Kastner, 2011; Kelly *et al.*, 2014]. In this suggestion, we attribute
30 awareness to other people as a means of modeling their attentional states, just as we
31 attribute it to ourselves to model our own attentional states.

32 For example, Bill pays attention to a hamburger in front of him. That mechanistic
33 process of attention leads to certain external signs on Bill such as his gaze direction,
34 facial expression, body language, and verbal cues. If you are observing Bill, then
35 based on a synthesis of those many cues you attribute awareness to him. You have an
36 internal model informing you that Bill is aware of the hamburger.

37 Arguably, your ability to attribute awareness to someone else is foundational to
38 all other social thinking. Maybe you think Bill is angry. It is difficult to attribute
39 anger to him unless you first understand that he is aware of the unpleasant thing that
40 ought to make him angry. You cannot predict that he will shout at you unless you
41 first understand that he is aware of you. Maybe you think Bill intends to reach out
42 and grasp something. You cannot make that attribution of intention unless you

1 understand that he is aware of the object to be grasped. Maybe you think that
2 someone else thinks that you think that he is lying to you. That complicated back and
3 forth of social cognition depends on understanding that the other person has such a
4 thing as awareness and is aware of you, of your likely thoughts, and of his own
5 thoughts. Social cognition makes no sense and has no foundation without the
6 underlying attribution of awareness.

7 It is not yet clear when animals evolved the ability to attribute awareness to each
8 other. Since many species of birds are highly social, perhaps birds can attribute
9 awareness to other birds [Thom and Clayton, 2013]. Certainly many mammals can,
10 including humans. The last common ancestor of birds and mammals lived approxi-
11 mately 350 MYA, and therefore a reasonable guess is that the social attribution of
12 awareness first appeared sometime before that — though of course it could have
13 evolved independently in both groups.

14 In this extension of the attention schema theory, awareness first evolved to help
15 control one's own attention, and then gradually expanded into another use that has
16 ended up defining us humans socially and culturally. It gave us our concept of mind
17 and allowed us to live immersed in a society of the minds of other people.

18 In the human brain, there is some evidence of overlap between the areas res-
19 ponsible for attributing awareness to others and the areas necessary for one's own
20 awareness. This overlap in function is particularly evident in an area of the cerebral
21 cortex called the temporo-parietal junction (TPJ), more or less just above the ears
22 and about an inch in. The TPJ has been a scientific puzzle because of an apparent
23 conflict between two competing lines of research. In one line of research, it is involved
24 in constructing models of other people's minds [e.g., Brunet *et al.*, 2000; Ciaramidaro
25 *et al.*, 2007; Fletcher *et al.*, 1995; Gallagher *et al.*, 2000; Goel *et al.*, 1995; Saxe and
26 Kanwisher, 2003; Saxe and Wexler, 2005; Vogeley *et al.*, 2001]. In another line of
27 research, the TPJ is involved in attention and awareness [e.g., Astafiev *et al.*, 2006;
28 Corbetta *et al.*, 2000; Mitchell, 2008; Shulman *et al.*, 2010]. Damage to the TPJ can
29 even cause a severe and long-lasting deficit in awareness called hemispatial neglect
30 [Karnath *et al.*, 2001; Vallar and Perani, 1986]. In neglect, typically damage to
31 the right side of the brain causes a loss of awareness of anything to the left side of
32 the body.

33 Why should a region of the cortex be involved in social cognition in some exper-
34 iments and in attention and awareness in other experiments? One possible reason
35 might be that this brain region participates in computations about awareness,
36 whether you are attributing awareness to yourself or to someone else. It would not be
37 correct to claim that the TPJ is the source of all computations related to awareness.
38 However, it may play a role.

39 We recently conducted an experiment to test this hypothesis more directly [Kelly
40 *et al.*, 2014]. The experiment involved two stages. First, people were scanned in an
41 MRI machine to measure brain activity. The subjects looked at a picture of a cartoon
42 face that was next to an object and rated how aware the cartoon person seemed to be

1 of the object. In this task, certain areas of the brain became active above control
2 levels. One area of activation was consistently within the TPJ.

3 In the second part of the experiment, the same people were taken out of the
4 scanner environment and placed in a different testing room. The hotspot in the
5 TPJ that was identified in the first part of the experiment was then targeted with
6 a technique called transcranial magnetic stimulation (TMS). In that technique, a
7 magnetic pulse is directed through the skull to temporarily disrupt brain function
8 in a small patch of tissue, approximately 1 cm wide. In this experiment, disrupting
9 the TPJ on one side of the brain disrupted the subject's ability to report dots
10 flashed on a screen on the other side of space. The effect was not general to the
11 entire TPJ. Instead, disruption of the specific hotspot obtained in the first part of
12 the experiment was necessary. When the disruption was targeted to another site,
13 2 cm away but still within the larger area of the TPJ, the effect was no longer
14 obtained.

15 One way to summarize this experiment is that specific areas of the brain became
16 active when a person looked at someone else and answered the question, "Is he aware
17 of the item next to him?" When the same brain regions were disrupted, the person
18 was less able to answer the question, "Am I aware of the item in front of me?" This
19 finding helps to support the hypothesis that awareness has taken on a social role at
20 least in humans. A system in the human brain participates in computations about
21 awareness whether you are attributing it to yourself or to someone else.

22 23 **8. Some Thoughts on Machine Consciousness**

24 In the attention schema theory, consciousness is more than a philosophical flourish. It
25 is a fundamental part of the data processing machinery of the brain. If the theory is
26 correct, then awareness is an internal model of attention and is crucial for the proper
27 regulation of attention. In addition, awareness has taken on ever-expanding roles
28 through evolutionary time including promoting the integration of information across
29 different domains and promoting social cognition.

30 All of these functions are as useful to artificial intelligence as they are to human
31 intelligence. They are also amenable to engineering. Every process described in this
32 paper could be built, though probably at first only at a simple level.

33 There is no fundamental or theoretical limit to stop computer scientists from
34 building a device that employs a human-like attention. In that process, signals
35 compete at a local and global level. Winning signals rise in strength and have a
36 disproportionate effect on memory and response choice.

37 There is also no fundamental or theoretical limit to stop engineers from adding
38 an attention schema to help that artificial device predict and therefore regulate its
39 own attention. That attention schema could contain simplifying information, mod-
40 eling attention as though it were an ectoplasmic and magical substance that can
41 reach out and "know" or "experience" things while being physically seated inside the
42 machine itself.

1 There is no theoretical limit to stop engineers from adding the equivalent of higher
2 cognition, a general purpose processor that is informed by the internal models
3 computed within deeper levels of the device.

4 Given these pieces, we would have a machine that is convinced it has subjective
5 awareness. If that higher cognition has access to language production, then the
6 machine could tell us that it has awareness. It would report that when it introspects it
7 finds awareness inside itself. It just knows it. Awareness is supplied to it *a priori*, like
8 a Kantian prior. It would behave, in these respects, like any person.

9 The device could be designed to attribute awareness not just to itself, but to others
10 as well. In that way, the machine would have a better basis for predicting the
11 behavior of others and also a more human-like social capability as it attributes spirit
12 to the beings around it.

13 The naïve approach of waiting to see if computers become conscious as they
14 become more complicated has not yet yielded a satisfactory result. It may be more
15 effective to design a machine in such a way that it concludes it has consciousness and
16 can report that conclusion. The machine could use that self-model to regulate its own
17 data flow and to understand the behavior of others.

18 If Deep Blue can beat Gary Kasparov, and Watson can win at Jeopardy, then a
19 computer that contains the essential components of consciousness is easily within
20 present technology. A concerted effort with sufficient resources could build such a
21 device, perhaps within a decade.

22 **Acknowledgment**

23 Supported by the Princeton Neuroscience Institute.

24 **References**

- 25 Ansong, U. and Heumann, M. [2006] “Shifts of visuospatial attention to invisible (metacon-
26 trast-masked) singletons: Clues from reaction times and event-related potentials,” *Adv.*
27 *Cogn. Psych.* **2**, 61–76.
- 28 Astafiev, S. V., Shulman, G. L. and Corbetta, M. [2006] “Visuospatial reorienting signals in the
29 human temporo-parietal junction are independent of response selection,” *European J.*
30 *Neurosci.* **23**, 591–596.
- 31 Baars, B. J. [1983] “Conscious contents provide the nervous system with coherent, global
32 information,” in *Consciousness and Self-Regulation*, eds. Davidson, R. J., Schwartz, G. E.
33 and Shapiro, D. (Plenum Press, NY), p. 41.
- 34 Baars, B. J. [1997] “Some essential differences between consciousness and attention, percep-
35 tion, and working memory,” *Conscious. Cogn.* **6**, 363–371.
- 36 Barlow Jr., R. B. and Fraioli, A. J. [1978] “Inhibition in the Limulus lateral eye *in situ*,”
37 *J. Gen. Physiol.* **71**, 699–720.
- 38 Beck, D. M. and Kastner, S. [2009] “Top-down and bottom-up mechanisms in biasing com-
39 petition in the human brain,” *Vis. Res.* **49**, 1154–1165.
- 40 Botvinick, M. and Cohen, J. [1998] “Rubber hands ‘feel’ touch that eyes see,” *Nature* **391**, 756.
- 41 Brunet, E., Sarfati, Y., Hardy-Baylé, M. C. and Decety, J. [2000] “A PET investigation of
42 the attribution of intentions with a nonverbal task,” *Neuroimage* **11**, 157–166.

- 1 Budd, G. E. [2008] “The earliest fossil record of the animals and its significance,” *Philos.*
2 *Trans. Roy. Soc. London B, Biol. Sci.* **363**, 1425–1434.
- 3 Ciaramidaro, A., Adenzato, M., Enrici, I., Erk, S., Pia, L., Bara, B. G. and Walter, H. [2007]
4 “The intentional network: How the brain reads varieties of intentions,” *Neuropsychologia*
5 **45**, 3105–3113.
- 6 Corbetta, M., Kincade, J. M., Ollinger, J. M., McAvoy, M. P. and Shulman, G. L. [2000]
7 “Voluntary orienting is dissociated from target detection in human posterior parietal
8 cortex,” *Nature Neurosci.* **3**, 292–297.
- 9 Crick, F. and Koch, C. [1990] “Toward a neurobiological theory of consciousness,” *Seminars*
10 *Neurosci.* **2**, 263–275.
- 11 Damasio, A. [1999] *The Feeling of What Happens: Body and Emotion in the Making of Con-*
12 *sciousness* (Harcourt, New York).
- 13 De Brigard, F. and Prinz, J. [2010] “Attention and consciousness,” *Wiley Interdiscip. Rev.*
14 *Cogn. Sci.* **1**, 51–59.
- 15 Dennett, D. C. [1992] *Consciousness Explained* (Back Bay Books, New York).
- 16 Desimone, R. and Duncan, J. [1995] “Neural mechanisms of selective visual attention,” *Ann.*
17 *Rev. Neurosci.* **18**, 193–222.
- 18 Engel, A. K. and Singer, W. [2001] “Temporal binding and the neural correlates of sensory
19 awareness,” *Trends Cogn. Sci.* **5**, 16–25.
- 20 Fletcher, P. C., Happé, F., Frith, U., Baker, S. C., Dolan, R. J., Frackowiak, R. S. and Frith,
21 C. D. [1995] “Other minds in the brain: A functional imaging study of ‘theory of mind’ in
22 story comprehension,” *Cognition* **57**, 109–128.
- 23 Franklin, G. F., Powell, J. D. and Workman, M. L. [1989] *Digital Control of Dynamic Systems*
24 (Addison-Wesley, Boston).
- 25 Gallagher, H. L., Happé, F., Brunswick, N., Fletcher, P. C., Frith, U. and Frith, C. D. [2000]
26 “Reading the mind in cartoons and stories: An fMRI study of ‘theory of mind’ in verbal and
27 nonverbal tasks,” *Neuropsychologia* **38**, 11–21.
- 28 Gazzaniga, M. S. [1970] *The Bisected Brain* (Appleton Century Crofts, New York).
- 29 Goel, V., Grafman, J., Sadato, N. and Hallett, M. [1995] “Modeling other minds,” *Neuroreport*
30 **6**, 1741–1746.
- 31 Graziano, M. S. A. [2013] *Consciousness and the Social Brain* (Oxford University Press,
32 New York).
- 33 Graziano, M. S. A. [2014] “Speculations on the evolution of awareness,” *J. Cogn. Neurosci.*
34 **26**, 1300–1304.
- 35 Graziano, M. S. A. and Botvinick, M. M. [2002] “How the brain represents the body: Insights
36 from neurophysiology and psychology,” in *Common Mechanisms in Perception and Action:*
37 *Attention and Performance XIX*, eds. Prinz, W. and Hommel, B. (Oxford University Press,
38 Oxford), pp. 136–157.
- 39 Graziano, M. S. A. and Kastner, S. [2011] “Human consciousness and its relationship to social
40 neuroscience: A novel hypothesis,” *Cogn. Neurosci.* **2**, 98–113.
- 41 Hwang, E. J. and Shadmehr, R. [2005] “Internal models of limb dynamics and the encoding of
42 limb state,” *J. Neural Engng.* **2**, S266–S278.
- Ivanoff, J. and Klein, R. M. [2003] “Orienting of attention without awareness is affected by
measurement-induced attentional control settings,” *J. Vis.* **3**, 32–40.
- Jacobs, O. L. R. [1993] *An Introduction to Control Theory* (Oxford University Press, Oxford).
- Jiang, Y., Costello, P., Fang, F., Huang, M. and He, S. [2006] “A gender- and sexual orien-
tation-dependent spatial attention effect of invisible images,” *Proc. Natl. Acad. Sci. USA*
103, 17,048–17,052.
- Karnath, H. O., Ferber, S. and Himmelbach, M. [2001] “Spatial awareness is a function of the
temporal not the posterior parietal lobe,” *Nature* **411**, 950–953.

- 1 Kawato, M. [1999] "Internal models for motor control and trajectory planning," *Curr. Opin.*
2 *Neurobiol.* **9**, 718–727.
- 3 Kelly, Y. T., Webb, T. W., Meier, J. D., Arcaro, M. J. and Graziano, M. S. A. [2014]
4 "Attributing awareness to oneself and to others," *Proc. Natl. Acad. Sci. USA* **111**,
5 5012–5017.
- 6 Kentridge, R. W., Heywood, C. A. and Weiskrantz, L. [1999] "Attention without awareness in
7 blindsight," *Proc. Biol. Sci.* **266**, 1805–1811.
- 8 Kentridge, R. W., Heywood, C. A. and Weiskrantz, L. [2004] "Spatial attention speeds dis-
9 crimination without awareness in blindsight," *Neuropsychologia* **42**, 831–835.
- 10 Koch, C. and Tsuchiya, N. [2007] "Attention and consciousness: Two distinct brain processes,"
11 *Trends Cogn. Sci.* **11**, 16–22.
- 12 Lackner, J. R. [1988] "Some proprioceptive influences on the perceptual representation of body
13 shape and orientation," *Brain* **111**, 281–297.
- 14 Lambert, A., Naikar, N., McLachlan, K. and Aitken, V. [1999] "A new component of visual
15 orienting: Implicit effects of visual information and subthreshold cues on covert attention,"
16 *J. Exp. Psychol. — Human Percep. Perform.* **25**, 321–340.
- 17 Lamme, V. A. [2003] "Why visual attention and awareness are different," *Trends Cogn. Sci.* **7**,
18 12–18.
- 19 Mack, A. and Rock, I. [1998] *Inattentional Blindness* (MIT Press, Cambridge).
- 20 McCormick, P. A. [1997] "Orienting attention without awareness," *J. Exp. Psychol. —*
21 *Human Percep. Perform.* **23**, 168–180.
- 22 Mele, S., Savazzi, S., Marzi, C. A. and Berlucchi, G. [2008] "Reaction time inhibition from
23 subliminal cues: Is it related to inhibition of return?" *Neuropsychologia* **46**, 810–819.
- 24 Merikle, P. M. and Joordens, S. [1997] "Parallels between perception without attention and
25 perception without awareness," *Conscious. Cogn.* **6**, 219–236.
- 26 Mitchell, L. P. [2008] "Activity in the right temporo-parietal junction is not selective for
27 theory-of-mind," *Cerebral Cortex* **18**, 262–271.
- 28 Mole, C. [2008] "Attention in the absence of consciousness?" *Trends Cogn. Sci.* **12**, 44–45.
- 29 Mulckhuysse, M. and Theeuwes, J. [2010] "Unconscious attentional orienting to exogenous
30 cues: A review of the literature," *Acta Psychol.* **134**, 299–309.
- 31 Mysore, S. P. and Knudsen, E. I. [2013] "A shared inhibitory circuit for both exogenous and
32 endogenous control of stimulus selection," *Nature Neurosci.* **16**, 473–478.
- 33 Newman, J. and Baars, B. J. [1993] "A neural attentional model for access to consciousness: A
34 global workspace perspective," *Concepts Neurosci.* **4**, 255–290.
- 35 Posner, M. I. [1994] "Attention: The mechanisms of consciousness," *Proc. Natl. Acad. Sci.*
36 *USA* **91**, 7398–7403.
- 37 Prinz, J. [2011] "Is attention necessary and sufficient for consciousness?" in *Attention: Phi-*
38 *losophical and Psychological Essays*, eds. Mole, C., Smithies, D. and Wu, W. (Oxford
39 University Press, Oxford), pp. 174–204.
- 40 Saxe, R. and Kanwisher, N. [2003] "People thinking about thinking people: fMRI investi-
41 gations of theory of mind," *Neuroimage* **19**, 1835–1842.
- 42 Saxe, R. and Wexler, A. [2005] "Making sense of another mind: The role of the right temporo-
parietal junction," *Neuropsychologia* **43**, 1391–1399.
- Scheidt, R. A., Conditt, M. A., Secco, E. L. and Mussa-Ivaldi, F. A. [2005] "Interaction of
visual and proprioceptive feedback during adaptation of human reaching movements,"
J. Neurophysiol. **93**, 3200–3213.
- Schiff, N. D. [2008] "Central thalamic contributions to arousal regulation and neurological
disorders of consciousness," *Annals NY Acad. Sci.* **1129**, 105–118.

14 *M. S. A. Graziano & T. W. Webb*

- 1 Shulman, G. L., Pope, D. L., Astafiev, S. V., McAvoy, M. P., Snyder, A. Z. and Corbetta, M.
2 [2010] “Right hemisphere dominance during spatial selective attention and target detection
3 occurs outside the dorsal frontoparietal network,” *J. Neurosci.* **30**, 3640–3651.
- 4 Thom, J. M. and Clayton, N. S. [2013] “Re-caching by Western scrub-jays (*Aphelocoma*
5 *californica*) cannot be attributed to stress,” *PLoS One* **8**(1), e52936.
- 6 Tononi, G. [2008] “Consciousness as integrated information: A provisional manifesto,” *Bio-*
7 *logic. Bullet.* **215**, 216–242.
- 8 Vallar, G. and Perani, D. [1986] “The anatomy of unilateral neglect after right-hemisphere
9 stroke lesions, a clinical/CT-scan correlation study in man,” *Neuropsychologia* **24**,
10 609–622.
- 11 van Boxtel, J. J., Tsuchiya, N. and Koch, C. [2010] “Consciousness and attention: On suffi-
12 ciency and necessity,” *Front. Psychol.* **1**, 217.
- 13 van Swinderen, B. [2012] “Competing visual flicker reveals attention-like rivalry in the fly
14 brain,” *Front. Integ. Neurosci.* **6**, 96.
- 15 Vogeley, K., Bussfeld, P., Newen, A., Herrmann, S., Happé, F., Falkai, P., Maier, W., Shah,
16 N. J., Fink, G. R. and Zilles, K. [2001] “Mind reading: Neural mechanisms of theory of mind
17 and self-perspective,” *Neuroimage* **14**, 170–181.
- 18 Wolpert, D. M., Ghahramani, Z. and Jordan, M. I. [1995] “An internal model for sensorimotor
19 integration,” *Science* **269**, 1880–1882.
- 20 Woodman, G. F. and Luck, S. J. [2003] “Dissociations among attention, perception, and
21 awareness during object-substitution masking,” *Psychol. Sci.* **14**, 605–611.
- 22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42