

Avshalom C. Elitzur

CONSCIOUSNESS MAKES A DIFFERENCE:
A RELUCTANT DUALIST'S CONFESSION

To Robert Jahn and Brenda Dunne

Introduction: Advancing the Mind-Body Problem into the Realm
of Science

If something odd persists, would its mere persistence make it natural? That would be the case with the layperson, but the scientist and philosopher should know better. Commonness should never mislead us to get used to the incredible.

Such is the phenomenon known as “consciousness,” underlying the age-old “Mind-Body Problem.” That consciousness exists at all is as odd today as it has been in ancient times. But here too, familiarity breeds contempt; the presence of consciousness at every moment in our waking lives often makes us forget how incredible it is.

For more than two millennia, the study of this problem (Block et al., 1997) has made no scientific progress. Physicalism,¹ dualism, and all other isms keep debating on it without being able to propose any decisive argument, not to mention experimental test, which could conclude the debate in favor of one theory or another.

But is this stalemate inevitable? I believe I have a scientific argument (Elitzur, 1989, 1996) in favor of one of the rival parties. Unfortunately, this party is interactionist dualism, which I dislike most. Indeed my argument comes with the expected penalty on this option, namely, entailing violation of a very basic physical principle. Being a physicist, this violation upsets me most.

Yet the argument is scientific, in that it derives, from a philosophical statements, an empirical prediction via the following reasoning:

1. By physicalism, consciousness and brain processes are identical.
2. Whence, then, the dualistic bafflement about their apparent nonidentity?
3. By physicalism, this nonidentity, and hence the resultant bafflement, must be due to error.

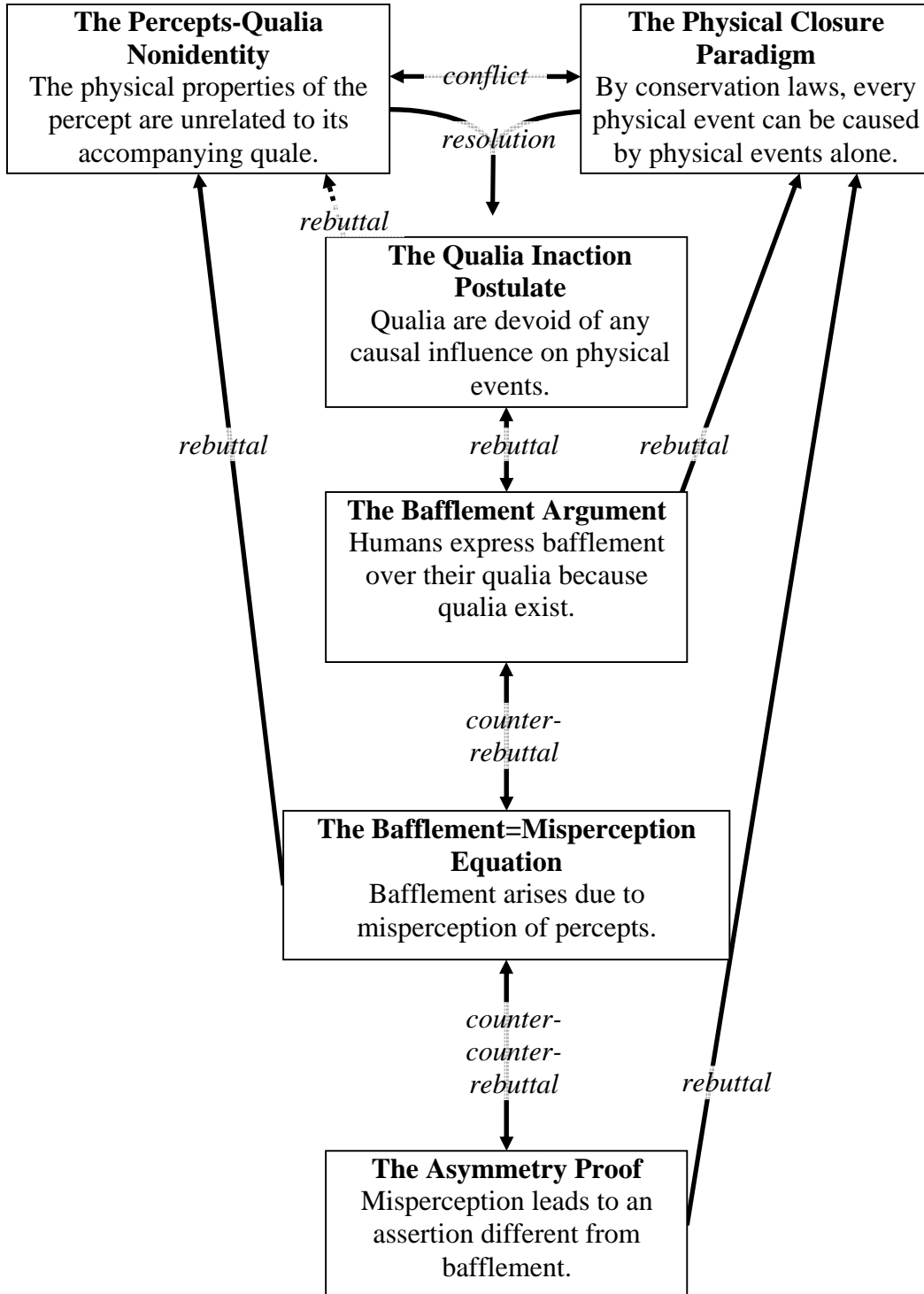
4. But then, again by physicalism, an error must have a causal explanation.
5. Logic, cognitive science and AI are advanced enough nowadays to provide such an explanation for the alleged error underlying dualism, and future neurophysiology must be able to point out its neural correlate.

This prediction, if rigorous, is falsifiable and therefore turns both physicalism and dualism into scientific theories in the full Popperian sense. Now, can this prediction further be falsified? I believe it can, even before the above disciplines respond to the challenge. This can be done with the aid of another powerful scientific procedure, namely, thought-experiment (Brown, 2007). Employing this procedure, I will show that no logical, cognitive or neural failure can produce the brain-consciousness nonidentity and the resulting bafflement as expressed by many humans. Ergo, bafflement about consciousness is a case where consciousness, as nonidentical with brain processes, exerts a causal effect of its own.

This paper's outline is as follows. In sections 1-3 I give an exposition of the Mind-Body Problem, with emphasis on what I believe to be the heart of the problem, namely, the Percepts-Qualia Nonidentity and its incompatibility with the Physical Closure Paradigm. In 4 I present the "Qualia Inaction Postulate" underlying all non-interactionist theories that seek to resolve the above problem. Against this convenient postulate I propose in section 5 the "Bafflement Argument," which is this paper's main thesis. Sections 6-11 critically discuss attempts to dismiss the Bafflement Argument by the "Bafflement=Misperception Equation." Section 12 offers a refutation of all such attempts in the form of a concise "Asymmetry Proof." Section 13 points out the bearing of the Bafflement Argument on the evolutionary role of consciousness while section 14 acknowledges the price that has to be paid for it in terms of basic physical principles. Section 15 summarizes the paper, pointing out the inescapability of interactionist dualism.

The scheme in Figure 1 gives this outline visually.

Figure 1



1. What's Your Mind-Body Problem Anyway?

Often, stating a problem well is half the way to its solution. Equally often, the Mind-Body Problem is ill-stated. Chalmers (1996), with typical wit, has shown that, when an author claims to have “solved” the Mind-Body Problem, they probably do not understand it in the first place. Chalmers then introduced his by-now classic distinction between the “hard problem” and the “easy problems” of consciousness. The “hard problem” is the one discussed below, whereas “easy problems” are

How does the brain process environmental stimulation? How does it integrate information? How do we produce reports on internal states? These are important questions, but to answer them is not to solve the hard problem: Why is all this processing accompanied by an experienced inner life? (pp. xi-xii).

Let us, then, present the “hard problem” first, so as to provide a basis for the arguments to come. I shall invoke a naïve discussant whose questions will help us focus on the crucial issues.

So what's the Mind-Body Problem with you anyway? Why don't you believe that science gives a satisfactory explanation of consciousness?

When dealing with consciousness, science miserably fails in what has always been its hallmark of success, namely, reducing qualities to quantities. “The qualitative” said Lord Rutherford, “is nothing but poor quantitative.” This, indeed, is usually the case. Consider, e.g., the following statements:

1. Red differs from blue.
2. Sweet differs from salty.
3. Love differs from hate.

These differences seem to be qualitative, but the scientific account neatly converts them into different numeric values on the same scales:

1. Both red and blue light are electromagnetic waves, differing only in their wavelengths: 700 nm for red and 400 nm for blue. Consequently, different cones in our retina react differently to these wavelengths due to different amino-acid sequences of their rhodopsin.
2. A sugar molecule, $C_6H_{12}O_6$, contains carbon, hydrogen and oxygen atoms, while a salt molecule, NaCl, contains sodium and chlorine atoms. All these atoms contain identical electrons on their shells, differing only in their numbers, which they exchange with the molecules in our tongue receptors.

3. Both love and hate involve very similar neurons, differing mainly in their location and spatial arrangement (location, specified by geometry, is also a quantitative measure).

In all these examples, qualitative differences between percepts turn out to be basically quantitative.

Thanks! I'll remember that next time I eat ice cream or hate someone. So why aren't you satisfied with the physical explanations to conscious experience?

While these explanations do a good job with percepts, rendering them (through neuroscience and chemistry) physical events, some intriguing phenomena that accompany these percepts are left out. These are pure qualities, qualia.

What's that? And what's the difference between qualia and percepts?

Qualia ("quale" in singular) are those aspects of our experience that cannot be communicated yet we know they are there. Suppose you and I look at a rose. Having verified that our color vision and linguistic abilities are normal, we assure each other that we both see a red rose. Still, you cannot rule out the possibility that I experience red the way you experience blue. True, in all languages each of us would name all colors the same way as the other. But this only means that we both have correctly learned to associate the appropriate word to the wavelength in question. Nothing of all that can tell you anything about my quale of the color. This is the notorious "inverted qualia problem."

The same holds for all other percepts of sound, smell, etc. The percept itself can be accurately communicated, but the accompanying quale remains inaccessible.

So, it's merely a problem of communication.

Much worse. Qualia elude not only communication, but observation and experiment as well. Suppose that, with sufficiently advanced technology, you obtain the fullest real-time description of what goes on in my brain – every neuron, synapse and neurotransmitter molecule – when I see a red rose. We have thus broadened the meaning of "percept" to the entire neurophysiological process that occurs when the stimulus is processed in the brain. Paradoxically, the problem now becomes worse:² You know better than I do what goes on in my brain when I perceive red, and yet, that doesn't bring you any closer to my quale of red.

Worse still, it is not only that you cannot be sure that my qualia are similar to yours – you cannot even be sure that I have any qualia at all.

With today's technology, a machine is perfectly conceivable that will name colors, in any language, with much greater accuracy than all humans. Does such a machine have the qualia of "red" or "blue"?

Returning to humans, the above "inverted qualia problem" leads to the even more grotesque "absent qualia problem," also known as "the problem of other minds." Personally I have no doubt that you, apart from appropriately responding to colors, sounds, tastes and odors, also experience their accompanying qualia. I am likewise sure that you feel happy when you laugh, besides the physical manifestation of laughter; that you are sad when crying, etc. And yet, even this very reasonable intuitive belief has no rigorous proof.

Isn't there a law that obliges qualia to come with percepts?

No, just as there is no law obliging qualia to come with thunderstorms or soap bubbles. Moreover, once you assume that the brain operates in compliance with physical law, qualia must not play any role in the brain's operation.

Here is why. Consider first the motions of billiard balls. Must you invoke any quale in order to explain them? Should you hypothesize that the balls "feel repulsion" upon colliding, or "yearn" to come to rest when slowing down? Their behavior is strictly and solely governed by the laws of mechanics. Next consider a plant that has not been watered for several days, nearly dying. You water it, and soon its leaves stretch again and regain their vitality. Should you invoke the qualia of "thirst" or "slaking thirst" to explain what happened? The physical laws governing osmosis (different concentrations of salt on the two sides of the cell's semi-permeable membrane) perfectly suffice.

You can guess where I am heading. Much higher up the scale of complexity, above balls and plants, are humans. Their percepts, no matter how complex, are supposed to be governed by neuronal processes that are, in essence, physical. Now you want to explain a certain behavior, say, picking a red rose. If your explanation invokes not only the percept of red but the accompanying quale as well, this amounts to asserting that the laws of physics do not sufficiently account for a physical process.

Why would that be so bad?

Well, if a non-physical cause plays a role in any process, than some of physics' most revered laws, such as energy and momentum conservation, are violated. It is easy to understand it in the case of the balls: If anything other than mechanical forces is involved with their motions,

then energy and momentum conservation must be violated. It would be much more difficult to prove such a violation with the plant drinking water, but the same violation must be involved in this case too.

Now let's return to the human picking a red rose. As long as only the percept of red affects her picking, then the accompanying quale plays no causal role and may be ignored. But if the quale too takes part, then the continuous, omnipresent causal network dominating physical reality must be somewhere broken. Somewhere along the neuronal chain, a physical event must occur that is not fully determined by the previous physical events. The very principle of causality is thereby violated.

But surely there is a difference between a few balls and a human! Our behavior is so complex...

Don't make the common error of letting qualia hide beneath complexity. The UN administration is many times more complex than each of its single officers. We may ascribe this administration various percepts – say, the UN “knows” and is even “concerned” about a war breaking – but it would be silly to ascribe the UN the quale of “concern” over such an event, even though a formal UN announcement may well express such concern. Complexity seems to be a necessary but not sufficient condition for qualia. Therefore it cannot explain why they exist.

True, the fact that we know qualia only from complex organisms like us makes the problem more delicate. The situation brings to mind G. B. Shaw's remark: “If you are not a communist at 20, you have no heart; if you are still a communist at 30, you have no head.” A similar choice awaits any self-consistent position with respect to qualia and levels of organization: Should you be silly or inhumane? If you grant qualia to, say, mice – believing that they have the quale of fear from cats, you may also ascribe the quale of “fear of light” to the photophobic response of a green alga, supposed to be accounted for by the automatic responses of its flagella.³ Or you may ascribe the quale of “fear of water” to a hydrophobic detergent molecule, supposed to be governed by electrical forces alone. On the other hand, if you deny the quale of fear of cats to mice, you may as well deny the quale of “fear of tigers” to a terrified Mogley running for his life.

Isn't the percepts-qualia distinction similar to the hardware-software distinction?

Certainly not! The shallow similarity between the two distinctions has misled many authors. Software, just like hardware, is a physical configuration of matter. All humans have software in the form of a brain within

which neurons are arranged and connected in special ways. Such a software is in principle observable to other persons. Furthermore, its existence is obliged by the laws of physics (biology, for this matter, being physics too). Qualia, on the other hand, remain unobservable and alien to physical law.

The problem, then, seems to be related to the logical problems involved with self-observation.

Again, I strongly disagree. All kinds of self-reference are prominent properties of human mind, and surely very interesting ones, but there is nothing paradoxical to them. Any intelligent system must refer to itself, yet self-observation and self-knowledge may or may not be accompanied with qualia. The former are examples of the “easy problems” while the latter present the “hard problem.” It is qualia which are strange, not entailed by any physical law...

But you keep talking about “physics” as if present-day physics is the final word. Can’t you imagine that future physics will reveal new phenomena, say, some unknown properties of matter or energy, which will eventually account for qualia?

You don’t have to be conversant in physics, nor in any scientific discipline, in order to realize that no physical concept, whether known or still to be discovered, can account for qualia. Imagine attending a lecture by a distinguished physicist. You know that physicists ascribe to particles some properties to which they give fancy names like “beauty,” “charm,” etc. So this physicist has discovered a new property, “loveliness,” and argues that this property, inherent also to the particles within the brain, accounts for our qualia. A frightful mathematics follows, explaining what “loveliness” is and how it gives rise to qualia. Naturally, if you have not specialized in particle physics, no chance you will understand what she is talking about. Or imagine that the lecturer is a world-renowned neurophysiologist announcing the discovery of a new neurotransmitter, say, alpha-mindo-encephaline, that explains qualia. Here too, a highly technical description follows, showing the complex pattern of interactions of that molecule with specific receptors within the neuronal synapses. And here too, unless you are a neurophysiologist, you will understand nothing.

Yet in both cases, if you look around in the audience, you will see senior scientists pensively nodding, perhaps raising some sophisticated objections, but finally saying, “Well, it’s interesting. Let’s think about it.”

Don't bother! Just ask, a priori: Can loveliness or alpha-whatever assure me, in principle, that my quale of red is not like your quale of blue? Worse, can it prove that any human has any qualia at all? Can any property of matter or energy rule out the possibility that my fellow humans lacks qualia altogether? You see, this is not a question of more knowledge.⁴ Qualia lie, in principle, beyond any possible physical account.

So why not ignore qualia altogether?

Well, let's see how long you can. Think about sleepwalking. Ridgway (1996) discusses in detail some cases where people committed murder while allegedly asleep, raising the question whether they can be held accountable for their deeds. Murder is an act requiring fairly advanced cognitive faculties, and yet, the people in these cases are believed to have been totally unaware of their own actions. Wilkes (1984) refers to a less substantiated case of a somnambulist physician who performed a medical examination and even made a correct diagnosis – all while asleep. In principle, there is no reason why this state cannot be extended to all mental functions. One might laugh, cry, – even converse! – while totally unconscious. For this reason, Wilkes titled her article “Is consciousness important?” Her answer being, expectedly, “No.”

Is she correct? Here is a thought experiment for you: How about turning all your qualia off, thereby putting you to permanent sleepwalking? All your percepts will remain the same, hence all your actions will be the same too. No one, therefore, would ever notice any difference in your behavior. But your qualia would be gone, forever. For doing this experiment on you, you will be paid \$1,000,000. Would you agree?

But it is doubtful whether it is possible to accurately...

Never mind technicalities! In physics, a *gedankenexperiment* (thought-experiment) is an indispensable tool, enabling one to anticipate technology by many years. Just bear in mind that physics allows the existence of humans with no qualia at all; and moreover that such humans are more compatible with physics than we conscious humans (see the discussion on zombies in Section 9 below). So, again, would you agree, for that nice sum, to go into lifelong sleepwalking which will leave your observed behavior intact but turn off your qualia of red and blue, sorrow and joy, forever?

Well, others won't see any difference, but for me it would be nothing short of death.

Welcome to the Mind-Body Problem! Science deals only with percepts, but it is qualia, which accompany every percept, that are the most essential ingredients of our life. Yet, they have no place even in the fullest and most detailed scientific explanation.

The Percepts-Qualia Nonidentity can be summarized as follows.

A Percept	A Quale
is a state occurring within one's brain upon perceiving something,	is the experience accompanying the percept,
obliged by physical law and evolving in strict compliance with it,	not entailed by physical law,
and can, in principle, be observed by another person, communicated and quantitatively measured with any desired accuracy.	and cannot, even in principle, be observed by another person, communicated or measured

Consciousness, then, is the totality of our qualia. Why is there consciousness at all, and how it is related to the brain, is the Mind-Body Problem.

2. No Room for Qualia: The Physical Closure Paradigm

But what is it that makes physics so inhospitable to consciousness? Among the most important pillars of physical law, and hence of all natural and life sciences, is a set of conservation laws, such as conservation of mass, energy, momentum and charge. It is due to these laws that some of the greatest scientific discoveries have been made. Conservation laws portray the universe as a closed system within which no mass, energy, etc. is ever added or subtracted. For this reason, every time some conserved quantity appears not to be conserved, a discovery may be on the way: Either a new phenomenon is to be revealed, or an important law is to be modified. This is the famous "closure of the physical world." Nothing other than physical phenomena – matter, energy, fields and their properties – takes place in the universal web of causes and effects. Because this assumption is so fundamental, underlying all natural and life sciences, I refer to it as not only a theory but rather as a paradigm.

Consider, then, a person responding to a stimulus. Anything other than physical causes affecting her response must lead to a violation of one of the conservation laws. Now, as that person also experiences

some qualia during the process, then, assuming the Percepts-Qualia Nonidentity, any causal role played by qualia is not only redundant but forbidden.

3. The Dilemma: Dismiss Qualia or Accept Violation of Physical Law

Between the Scylla of the Percepts-Qualia Nonidentity and the Charybdis of The Physical Closure Paradigm, lies the entire strait of the Mind-Body Problem. And there are Sirens too, in the form of several theories offering their solutions. They can be grouped into two major types, namely, physicalism and dualism.⁵ Physicalism invented a variety of exercises in order to prove that qualia do not really exist, being merely some aspects of the percepts (Dennett, 1991, 2003). Dualism, on the other hand, straightforwardly acknowledged that qualia exist alongside percepts.

Physicalism never won the full acceptance of the scientific and philosophical communities. If there is something to “red” that I cannot communicate and that even the most detailed description of my brain cannot yield – then something probably exists that lies outside of the framework of present-day science. Now these qualia, for any person, are by no means trivial: To lack them is to be inwardly dead. It would therefore be silly to ignore their existence.

While dualism does not dismiss qualia, the cure that it offers seems to be worse than the disease. It flies, as shown in the previous section, in the face of the Physical Closure Paradigm.

Understandably, some people turned to parapsychology in search of a straightforward proof that something non-physical interferes with behavior. Most notable were the Princeton Engineering Anomalies Research Laboratory (Jahn & Dunne, 1987, 2001). Yet, so far, after many years of admirable labor, the effects they found have not been strong enough to convince the scientific community.

Others turned to quantum mechanics for help. QM, so it seems, has undermined determinism, hence an interference of qualia with the brain’s random macroscopic events may not violate physical law after all. Thus Eccles (1994) has invoked QM to allow free will to interfere with the neurons’ synapses without violating energy conservation. Other quantum-mechanical models were proposed by eminent physicists (e.g., Penrose, 1994; Stapp, 2009; Tuszynski, 2006). Although being ingenious, they do not purport to resolve the “hard problem” (see the proverbial particle physicist in Section 1), but only argue that the action of

consciousness can preserve energy conservation, in compliance with the first law of thermodynamics.

The trouble, however, is that there is also a second law of thermodynamics, dealing with entropy increase. One of this law's derivatives, associated with the famous "Maxwell's demon" paradox (Fanchon et al., 2009), says that it is impossible to introduce order into a disordered process without investing energy. This leaves dualism with two options. Either

- a. Qualia's effect on behavior is random. That won't help much. To affect behavior, qualia must do so consistently. For example, if the quale of red, in addition to the percept, affects one's picking a red rose, then that quale must work the same way every time the percept appears.

Or

- b. Qualia's effect is systematic. But then, qualia must be using energy in order to interfere with the brain's random processes in a nonrandom manner, again violating the first and/or the second laws of thermodynamics.

To summarize, physicalism advises us to believe that qualia do not really exist, while dualism acknowledges their existence at the price of allowing physical anomalies.

4. The Consensual Compromise: The Qualia Inaction Postulate

So, is physicalism or dualism offering the lesser evil? To better address this question, we should consider a few variants/hybrids of these rival schools:

- a. Identity/double-aspect theory: The quale and the percept are one and the same process, only perceived as different.
- b. Parallelism: The quale and the percept are different processes, belonging to the mental and the physical realms, respectively. In each realm, events follow one another in a strict cause-and-effect manner. However, the two realms run parallel, never interfering with one another. Only by virtue of their perfect correlation they give rise to the illusion that they causally affect one another.
- c. Epiphenomenalism: The quale and the percept are different processes, belonging to the mental and the physical realms, respectively. But they maintain asymmetric causal relations: percepts give rise to qualia but never vice versa.

Now, we do not need to go into the details of these theories, for they all share with physicalism one crucial postulate: Qualia play no causal role. Let us examine this postulate with the aid of an ordinary behavior: A woman kisses a man. First consider the man-in-the-street explanation of this behavior:

a. Alice kisses Bob because she loves him (common sense).

The Mind-Body Problem arises with its full acuity: “Love” is both a percept and a quale. Which of them, then, constitutes the kiss’s cause? Let’s first pursue the most daring option:

b. Alice’s kissing Bob is caused by the quale of her love to him (Interactionist Dualism).

No scientifically minded scholar would accept such an account, fearing the clash with the Physical Closure Paradigm. To avoid this, one might try:

c. Alice kisses Bob because the percept of love, caused by sensory signals coming from him, triggers the behavior of kissing (Physicalism).

But then where, in all this, is love’s quale? A prudent theorist might clarify this position as follows:

d. Alice’s quale of loving is the percept of loving; she only perceives the two things as distinct (Identity or Double-Aspect Theory).

Or

e. Alice has the quale of loving alongside with the percept. It is only the percept, being a physical event, which gives rise to the consequent physical act of kissing, while the accompanying quale gives rise only to consequent qualia (Parallelism).

Or

f. Alice has the quale of loving alongside with the percept. It is only the percept, being a physical event, which gives rise to the consequent physical act of kissing, yet this percept, as well as all the consequent brain states, produce the accompanying qualia (Epiphenomenalism).

Options (e-f) are semi-dualist or noninteractionist-dualist theories. Their bottom line, shared with physicalism (c-d) as well, is this: For any instance where a quale seems to affect behavior, it can be shown that it is not the quale but its physical parallel, the percept, which exerts the effect. (Kim, 1996).

This can be succinctly put as the Qualia Inaction Postulate:

Any behavior would be exactly the same had there been no qualia.

Little wonder that this postulate has been opted for by most modern philosophers.⁶ It acknowledges that something very peculiar is going on within us, yet assures us that this thing has no influence whatsoever on our observable behavior. Consciousness, then, should make no difference.

5. But Qualia Do Play a Causal Role: The Bafflement Argument

The stakes are now very high. We need only one example in which the Qualia Inaction Postulate fails – where a quale alone, in itself, not its parallel percept, exerts a causal effect – to falsify all the comfortable alternatives to interactionist dualism. Guess what: This example is occurring to you, dear reader, at this very moment!

For, why do we feel and say that qualia are not identical with brain processes? Why do we talk, argue and write about the Mind-Body Problem? I submit that the answer is simply this: We are baffled by qualia because we have qualia. Hence, as against the Qualia Inaction Postulate, I propose the Bafflement Argument:

The fact that humans are baffled by the Percepts-Qualia Nonidentity, and express this bafflement by their observable behavior, is a case where qualia per se – as nonidentical with percepts – play a causal role in a physical process.

6. Let's Explain it Away: The Bafflement=Misperception Equation

Of course, adherents of the Qualia-Inaction Postulate and the Physical Closure Paradigm will not let us get away with it so easily. Anticipating their objections will help us advance the problem further into the scientific realm.

“He who increases knowledge increases pain” (Ecclesiastes 1, 18), but the opposite is equally correct. Isn't it significant, now that we come to think of it, that most discussions of qualia take a painful experience as a starting point? Happy experiences are taken for granted! So, for knowledge's sake, let us inflict the following pain on our Alice. She and Bob have broken up, leaving Alice sad over the separation. Now Alice wonders: Why is there a quale of sadness?

No matter how much neuroscience she studies, she only finds herself delving deeper into the Mind-Body Problem. In fact she has already found consolation in a new love, but the problem keeps baffling her,

hence she talks, argues and attends conferences on it. She may even write a paper.

Can we subject Alice's bafflement to the same procedures to which we earlier subjected her act of kissing? The man-in-the-street account is, again, simple:

a. Alice says that qualia are baffling because she experiences a quale that is nonidentical with her percept (Interactionist dualism).

This, again, is anathema to the Qualia Inaction Postulate, and hence to the entire Physical Closure Paradigm. There seems to be only one way to avoid such a clash with physical law, namely, to prove that Alice's bafflement is misguided, and, moreover, that the Percepts-Qualia Nonidentity is false. Consider the less interesting case in which Alice believes in evil eye. Surely no one would ascribe this belief to the existence of evil eye. Rather, we would ascribe it to an error, stemming from faulty reasoning. Why not, then, ascribe the bafflement about qualia to a similar error? Here, then, is the alternative:

b. Alice says that qualia are baffling because faulty reasoning misleads her to believe that qualia are nonidentical with percepts (Physicalism).

This position can be generalized to the Bafflement=Misperception Equation:

Some people misperceive their own percepts, falsely believing that they are accompanied with nonidentical qualia, therefore expressing bafflement over this duality.

7. Misperception Entails a Neural Failure: A Testable Prediction Derived from the Bafflement=Misperception Equation

Notice, first, that by this explanation-away the physicalist position commits itself, for the first time, to a falsifiable prediction: When future neurophysiology becomes advanced enough to point out the neural correlates of false beliefs, a specific correlate of this kind would be found to underlie the bafflement about qualia. Then, at long last, we shall understand why, for more than two millennia, numerous otherwise-intelligent authors misperceived their percepts, thereby being misled to believe in the Qualia-Percepts Nonidentity.⁷

I don't believe any of this. But the point is that the question is not matter of belief anymore. It has rather turned into an empirical issue. Conversely, the Bafflement Argument entails the opposite prediction: No neural pattern underlying a false belief will be found to underlie adherence to dualism.

Happily, we don't have to wait for future science to rule in this matter. Simple logic suffices to show that those who ascribe dualism to misperception commit an embarrassing error.

8. The Crucial Issue: If the Bafflement=Misperception Equation is Valid, the Mind-Body Problem is Over

It is now time to give the Bafflement=Misperception Equation its most due serious consideration: Is the expression of bafflement over qualia obliged by some physical laws governing our brains? Or, put in the AI language, Is there a physical/logical principle that obliges an intelligent system to express bafflement of this kind? I italicize these questions because they are of utmost importance; please bear with the didactic tone of my italicizing their significance too:

If a proof is ever given that an intelligent system, by virtue of physical laws alone, must state that it has qualia which are nonidentical with percepts, then the age-old Mind-Body Problem would finally get a definite solution – a physicalist one. The Percept-Qualia Nonidentity would turn out to be nothing but an unfortunate misperception, inherent to all intelligent systems, and the problem would turn out to be a pseudo-problem.

Yes, it will be that simple: Just as there is no “rabbit-duck problem,” “left-right positioned Necker cube problem” or any other problem entailed by misperception, so would the Mind-Body Problem finally turn out to be the mere failure of many otherwise-intelligent people to realize that percepts and qualia are just one and the same thing. All dualistic arguments made over the millennia – from Plato to Descartes to Leibniz to Popper and all others – would be put to rest by a simple and collective *ad hominem*, formulated in the precise terms of cognitive science. Dennett will triumphantly say: “I told you so!” Philosophy and science would then move on – for good – to other issues.

Now, the authors quoted in the following sections are alleging just such a proof, evidently without realizing that, if correct, this proof amounts to no less than The ultimate solution to the Mind-Body Problem. The only question is, Is this proof valid?

Don't hold your breath.

9. Chalmers: Zombies Misperceive Just as We Do

In his delightful *The Conscious Mind* (1996) Chalmers' briefly objects to my Bafflement Argument:

Indeed, Elitzur (1989) argues directly from the existence of claims about consciousness to the conclusion that the laws of physics cannot be complete, and that consciousness plays an active role in directing physical processes (he

suggests that the second law of thermodynamics might be false). But I have already argued that interactionist dualism is of little help in avoiding the problem of explanatory irrelevance (p. 183).

Chalmers has struggled with a similar idea in his discussion of “zombies.” These creatures are very instructive. Imagine intelligent beings that resemble us in every detail of our physiology, neuroanatomy and chemistry, but have no qualia.⁸ This, recall, is perfectly consistent with physics – in fact, as noted above, zombies accord with physics more than the existence of non-zombies.

Our question now assumes a specific form: Would zombies be as baffled by qualia as humans are?

Astonishingly, Chalmers’ answer is in the affirmative. His reasoning is so peculiar that I prefer to use lengthy quotes:

To see the problem in a particularly vivid way, think of my zombie twin in the universe next door. He talks about conscious experience all the time – in fact, he seems obsessed by it. He spends ridiculous amounts of time hunched over a computer, writing chapter after chapter on the mysteries of sensory qualia, professing a particular love of deep greens and purples. He frequently gets into arguments with zombie physicalists, arguing that their position cannot do justice to the realities of conscious experience.

And yet he has no conscious experience at all! In his universe, the physicalists are right and he is wrong. Most of his claims about conscious experience are utterly false. But there is certainly a physical or functional explanation of why he makes the claims he makes. After all, his universe is fully law-governed, and no events therein are miraculous, so there must be some explanation of his claims. But such explanations must ultimately be in terms of physical processes and laws, for these are the only processes and laws in his universe (p. 180).

Absurdity culminates with the conclusion:

The explanation of his claims obviously does not depend on the existence of consciousness, as there is no consciousness in his world. It follows that the explanation of my claims is also independent of the existence of consciousness (p. 180).

One has to read this passage time and again in order to believe what it says: A philosopher writes a book about qualia, discussing their enigmatic nature in great detail, and then states that he would write exactly the same book had he lacked qualia!⁹

But why on Earth should zombies express bafflement about qualia if they don’t have any? After all, Chalmers professes physicalism, by which there is a cause for anything zombies say. If the cause of the talk about qualia is not qualia themselves, what is it? As he himself frankly asks,

To get some feel for the situation, imagine that we have created computational intelligence in the form of an autonomous agent that perceives its environment and has the capacity to reflect rationally on what it perceives. What would such a system be like? Would it have any concept of consciousness, or any related notions?

His answer is “yes,” reasoning, in essence, that the zombie misperceives his “direct, unmediated” percept as distinct from what he knows about that percept:

[I]t seems likely that such a system would have the same kind of attitude toward its perceptual contents as we do toward ours, with its knowledge of them being direct and unmediated, at least as far as the system is concerned. When we ask how it knows that it sees the red tricycle, an efficiently designed system would say, “I just see it!” When we ask how it knows that the tricycle is red, it would say the same sort of thing that we would do: “It just looks red.” If such a system were reflective, it might start wondering about how is it that things look red, and about why it is that red just is a particular way, and blue another. From the system’s point of view it is just a brute fact that red looks one way, and blue another. Of course from our vantage point we know that this is just because red throws the system into one state, and blue throws it into another; but from the machine’s point of view this does not help (p. 185).

Let us summarize this reasoning:

1. People have qualia.
2. People express bafflement about qualia.
3. Physics allows the existence of zombies that lack qualia.
4. Zombies must also express bafflement about qualia.
5. Therefore people express bafflement about qualia for reasons other than their having qualia.

Chalmers is well aware (personal communication, July 2004) that this position is awkward. But the situation is worse. It is obvious that Chalmers’ zombie, by Chalmers’ own typology (see section 1 above), is addressing one the “easy problems” while the real Chalmers addresses the “hard problem.” This gives rise to a clear asymmetry: Just ask yourself: Can a zombie brood over the “problem of absent immediate perceptions” as we brood over the “problem of absent qualia”? By Chalmers’ own account, the answer is negative! This asymmetry will shortly enable us to prove that the Bafflement=Misperception Equation is plainly wrong.

10. Chalmers vs. Charmless: A Revised Turing Test

Notice, first, that Chalmers phrased his Zombie Universe too loosely. He did not even consider the possibility that the zombie philosopher will not express bafflement. Indeed, there are two definitions for a phenomenal zombie, namely,

- a. A human that behave exactly like an ordinary human but lacks qualia.
- b. A human that physically operates like an ordinary human but lacks qualia.

The difference is important. The first definition is opaque to physical reasoning, inevitably leading to the “zombie bafflement” mess discussed in the previous section. The second definition is more scientific in that it only specifies the system’s initial conditions and leaves the resulting state to be causally deduced from the former. Let us, therefore use the second definition.

We are now in a position to propose an experiment that, though not yet feasible, may transfer the study of qualia from philosophy to empirical science. In essence, it is a variant of the celebrated Turing test (Turing 1950), designed for judging whether a sufficiently advanced computer can simulate human intelligence. With the proposed modification, this test can give a clear-cut answer to a much more fundamental question.

Turing’s test was this: Let a computer and a human dwell in two separate rooms. Let the experimenter, unaware of their locations, send whatever questions she has in mind to both rooms via electric cables and get their answers in the same way. If she fails to tell by the answers who is the machine and who is the human, then the machine, for all practical purposes, is as intelligent as a human.

Suppose, now, that such a future computer has passed the test. The time is now ripe for the greatest question of all: Does this computer have qualia?

First, let us give our computer a name. How about Charmless? The slight difference from its human namesake indicates that, unless we prove that it has qualia, then, even if it is as bright and witty as Chalmers, it is Chalmers’ zombie incarnate.

The test is straightforward – a seemingly innocuous question, say, “What is red?”

Assuming that we have followed Turing's recommendation "to provide the machine with the best sense organs that money could buy" (Hodges, 1988), Charmless might give a very accurate answer, such as

1. "Red is the color I perceive when electromagnetic radiation of wavelength 700 nm impinges on the photoelectric device at the back of my obscure chamber, absorbed by photochemical molecules sensitive to this wavelength, and converted into electric pulses that go through optic fibers to the color-recognition system that in turn activates my memory, language and vocal systems."

This answer is much more detailed and accurate than that of an average human. But it would indicate that Charmless is a zombie, indeed devoid of phenomenal charm.

On the other hand, the above answer might have a baffled addendum, such as

2. "Red is... [see 1 above]. However, there is something to my subjective experience of red that is not indicated by the description I just gave you. I know of no way of communicating that additional ingredient. Although I can see that you and I refer to the same color when we use the same word, I can never be sure whether your subjective experience of red is not what I experience as blue. In fact, I am not sure you have any subjective experience accompanying your color perception."

This, as with humans, would indicate that Charmless too has qualia in addition to his percepts.

Notice that in this case we can determine with certainty the cause of this bafflement. Since Charmless is man-made, we can rule out the possibility that his bafflement is the result of some pre-installed "bug" such as an explicit command to express bafflement or some deliberate misperception imposed on it. In other words, we can rule out any cause to Charmless' assertion about having qualia other than his really having them.

11. The Crucial Question Sharpened: Would Zombies be Baffled Too?

Following Chalmers' analysis, a heated debate ensued over the issue of zombies, particularly over the question whether they might be baffled over qualia. Moody (1994), contra Chalmers, argued that they would not. Several articles followed Moody's, most of them expectedly objecting to his conclusion. Flanagan and Polger (1995) argued that zombies might wonder, just as conscious humans, whether qualia are inverted:

Suppose that normal zombies, upon seeing light of a certain wavelength x go into a state that is the disposition to say, "that object is green", and then they

act on that disposition. [...] All that is necessary for an inverted color judgment problem is that behavioral pathways get crossed twice. In our case (i.e., the usual inverted spectrum problem) one of the pathways is supposed to be the qualitative look of color, and the other a speech act. Zombies could have an equivalent problem with two non-conscious inversions [...]. First, when seeing an object that reflects a wavelength x , the inverted color judgment zombie enters the state that, in normal zombies, is the disposition to say, "that object is red." However, due to the second crossed wire, the inverted color judgment zombie's "that object is red" state actually causes it to utter, "that object is green." Thus a double inversion can create a problem indistinguishable from the inverted spectrum problem.

This primitive version of the Bafflement=Misperception Equation can be straightforwardly refuted: Just uncross the damn wires! Then let the zombies inspect their brains to see that no wires are crossed anymore, and they will stop worrying about inverted qualia for good. Now, do Flanagan and Polger believe that we humans have a similar cross-wiring? In this case neurophysiology is on its way to one of its greatest discoveries, and the Mind-Body Problem will shortly find its solution. I will not blame you, dear order, if you are not willing to bet on this possibility.

Dennett's (1995) attack on Moody was sharper, launched at the logical level:

If, *ex hypothesi*, zombies are behaviorally indistinguishable from us normal folk, then they are really behaviorally indistinguishable! They say just what we say, they understand what they say (or, not to beg any questions, they understand^z what they say), they believe^z what we believe, right down to having beliefs^z that perfectly mirror all our beliefs about inverted spectra, "qualia," and every other possible topic of human reflection and conversation (p. 322, italics original).¹⁰

Dennett correctly points out a flaw in Moody's formulation of the problem (see the two definitions for a zombie in section 10 above), but he still misses the essential point. Moody's question can be easily re-phrased so as to be immune to Dennett's criticism:

"Suppose that there are zombies that behave just as we do yet lack qualia. Would their bafflement about qualia be consistent with physical laws?"

Thus rephrased, the question has two possible answers, each of which having far-reaching consequences:

- a. Zombies will be baffled over qualia by virtue of some physical cause. In this case, bafflement about qualia is not due to the existence of qualia. But then, another cause for this allegedly erroneous bafflement must be

detected by future neurophysiology, which will indicate that we somehow misperceive our percepts. This is the result obliged by physicalism.

b. Zombies will not be baffled over qualia. Dualism would then be vindicated.

12. The Asymmetry Proof: Genuine Bafflement Has No Physical Cause

Thankfully, we do not need zombies to exist in order to give the cardinal question about zombie bafflement a definite answer. We are now in a position to prove that all the physicalist counter-arguments to the Bafflement Argument, such as the Bafflement–Misperception Equation, run into contradiction. Let us take Chalmers' argument as a representative example. Recall that by physicalism, the laws of perception are, in essence, physical laws too. Here, then, is the proof:

1. A presumably conscious human (henceforth Chalmers) states that his percept P is not identical with the corresponding quale Q.
2. Chalmers further argues that a zombie duplicate of him (henceforth Charmless) is possible, who has only P without Q.
3. Chalmers asserts that, by physical law, Charmless must notice a difference between what he knows about the physical process underlying his percept and the unmediated percept itself.
4. Chalmers then argues that this difference (3) must produce in Charmless the same bafflement as Chalmers' bafflement about the P-Q nonidentity (1).
5. Ask now Chalmers: Can you conceive of a Charmless who is identical to you but lack Q? His answer, by (2), is "Yes."
6. Next ask Charmless: Can you conceive of a duplicate of you (henceforth Harmless) who is identical to you but lacks Q? His answer, by (3), must be "No; unmediated percepts must occur by physical law."
7. As Chalmers can conceive of Charmless but Charmless cannot conceive of Harmless¹¹, the two kinds of bafflement, associated with (1) and (3), are essentially different.
8. Hence, the physical explanation for (3) does not hold for (1).
9. Hence, (4) is false.

The same contradiction can be shown to follow any theory that denies causal role to qualia. Let me reiterate the reasoning presented in this paper's Introduction: Why do we perceive qualia as distinct from

our neural firings? If it is not because they are different, then we have a misperception of them as different. But misperception, just like perception, may or may not have an accompanying quale. Take a simple optical illusion, say, a straight line appearing to be curved. Both the correct and the incorrect percepts may or may not have a quale. Now, the optical illusion is obliged by the laws of perception, whereas qualia are not!

Generalizing, we get the Asymmetry Proof:

If a quale is identical with its percept, then its appearance as nonidentical must be due to misperception. But misperception, being a special kind of perception, occurs in accordance with physical law. Hence, upon reflection, it must turn out to be obligatory. Qualia, in contrast, can be conceived of as altogether absent.

13. A Welcome Consequence: Qualia Make Evolutionary Sense

Before concluding, let us appreciate one of the Bafflement Argument's scientific benefits, which for upholders of the Qualia Inaction Postulate is unattainable.

If Alice kisses Bob only by virtue of neural mechanisms developed during evolution, whence the quale of loving? The question is well known in the broader biological context. If a rabbit escapes a fox only by virtue of neural mechanisms, and if the fox chases it unaffected by the quale of hunger, why are there qualia of fear and hunger in the first place?

In his famous "What is it like to be a bat?" Nagel (1974) has extended the question of qualia to animals, whose expressions of emotions seem to indicate that they have qualia too. What qualia – if any – has a bat when perceiving something with the aid of echolocation? An eminent zoologist and pioneer of bat echolocation research, Donald Griffin, responded to Nagel's challenge with some fascinating books (e.g., 1992) claiming to prove animal consciousness. However, most philosophers, even those sympathetic to this goal, pointed out that Griffin ignored the more serious "problem of other minds" (see section 1): There is no proof yet for the so natural assumption that other humans have consciousness!

Once, however, bafflement indicates the ability of qualia to affect behavior, then the question gets a very reasonable answer. Alice's kiss may take a bit longer thanks to the additional effect of love's quale, and the qualia of hunger and fear may add some speed to the rabbit's and the fox's race. Evolution is capable of magnifying even the minutest effects (Elitzur, 1994). Qualia, therefore, may give a real advantage for survival.

14. The Penalty: Conversation Laws Might Be Wrong

While interactionist dualism makes evolutionary sense, its clash with conservation laws, so fundamental to physics, is an unavoidable price, and a heavy one. All I can say in reply is that as far as the reasoning in this article was sound, no consequence should be feared. In addition, let me point out that physics itself is beset by several paradoxes that indicate a future revolution. Basic notions like space, time and causality must be thoroughly revised in order to resolve some of modern physics' longstanding paradoxes (Elitzur et al., 2005).

15. Facing the Inescapable: The Physical Explanation of Behavior is Incomplete

At the end of the day, it is astonishing that, throughout the enormous literature on the Mind-Body Problem, nearly no attention has been given to the so simple question: Why are people baffled over the Mind-Body Problem?

I submit that there are only two consistent answers. Either

- a. People are baffled over the Mind-Body Problem because the problem is genuine, i.e., qualia are not identical with percepts. Ergo, it is qualia's very existence that gives rise to the behaviors associated with bafflement. Ergo, something non-physical interferes with physical processes;

Or

- b. People are baffled over the Mind-Body Problem for erroneous reasons. Qualia are identical with the percepts. Ergo, the physical causes for the misperception of qualia will eventually be found, explaining once and for all the widespread failure to accept the percepts-qualia identity.

Answer (b), I submit, has been shown in this article to be flawed.

Then there is a chimera of (a) and (b), endorsed by Chalmers (1996) and Flanagan and Polger (1995):

- c. People are baffled over the Mind-Body Problem, and the problem is genuine, i.e., qualia are not identical with the percepts. However, people are baffled by this nonidentity for reasons other than this nonidentity.

Allow me to formulate (c) in common language:

"I have qualia, but I would have said that I had qualia even if I had not. Still, I have qualia. Believe me, I do!"

I must confess that have I heard such a statement from a person in the street I would suspect that he is schizophrenic or at least severely schizoid. That such a position is endorsed by competent philosophers only attests to the acuity of the cardinal yet neglected question as to why people are baffled by qualia. Dualism, alas, offers the most reasonable answer.

Notes

1. Previously known as “materialism.”
2. This is in essence Jackson’s “knowledge argument,” forcefully presented in his paper “What Mary Didn’t Know” (1986), which gave rise to an entire volume titled *There’s Something About Mary* (Ludlow et al., 2004).
3. Facilitated – how inventive is evolution! – by rhodopsin, the same pigments facilitating color vision in our eyes.
4. I am indebted to Uzi Awret for the following keen observation: “While the easy problem of consciousness is a result of knowing too little, the hard problem results from knowing too much.”
5. A more accurate dichotomy will be the contrast between physicalism (“the physical world is everything; mind is an illusion”) and mentalism (“everything is mind; matter is an illusion”). Conversely, one may distinguish between monism (“there is only one reality, namely, the physical/mental”) with dualism (“there are two realities, the physical and the mental”). But for our purpose the above dichotomy suffices.
6. Kim (1996) refers to this postulate as the “exclusion argument,” and Flanagan (1997) as “conscious inessentialism.”
7. A possible objection to this prediction can argue that a false belief can exist in one’s brain while its causes, namely, the previous events that have lead to this belief, no longer exist in that person’s memory. However, the experience of both psychodynamic and cognitive therapy shows that this is not so. The causation of false belief seems to be accessible to the therapist, albeit with considerable effort, their elucidation leading to the removal of that false belief. Neurophysiology will therefore follow suit.
8. Chalmers distinguishes between a psychological zombie and a phenomenal one. The zombie known from voodoo horror stories (or from the common derogatory term) is a psychological zombie, manifesting a clear behavior of a “living dead” such as apathy and lack of emotions. The zombie with which we deal, in contrast, is a

phenomenal zombie, capable of manifesting all emotions manifested by humans, while only lacking the respective qualia.

9. And to make the irony perfect, it is Chalmers who has added a mischievous comment in his book (p. 190) wondering whether the staunch physicalist Dennett is a zombie!

10. The superscript z denotes, following Chalmers, zombie mental functions that resemble ours but lack qualia.

11. Which is why we need not worry about Armless and so on.

Acknowledgements

It is a pleasure to thank the PEARlab members for their stimulating Conscious Academy on July 2003 at Princeton University, where the first draft of this article has been written, and to Harald Atmanspacher for a highly inspiring dialogue. The imaginary discussant invoked in Section 1 has been, in fact, Moran Cerf. Special thanks are due to Andreas Lindblom, Nili Alon and Vijay Chandrasekaran for their helpful comments. When this article was under its last revision I found out that Uziel Awret has reached similar conclusion in an unpublished manuscript. I thank him for enlightening discussions.

References

- Allen, C. (2006) Animal consciousness. Stanford Encyclopedia of Philosophy, <http://plato.stanford.edu/entries/consciousness-animal/>.
- Block, N., Flanagan, O., & Güzeldere [Eds.] (1997) *The Nature of Consciousness: Philosophical Debates*. Cambridge, Mass.: MIT Press.
- Brown, J. R. (2007) Thought experiments. Stanford Encyclopedia of Philosophy, <http://plato.stanford.edu/entries/thought-experiment/#SomRecVieThoExp>.
- Chalmers, D (1996) *The Conscious Mind*. Oxford: Oxford University Press.
- Dennett, D. C. (1991) *Consciousness Explained*. New York: Little, Brown.
- Dennett, D. C. (1995) The unimagined preposterousness of zombies: Commentary on T. Moody, O. Flanagan and T. Polger. *Journal of Consciousness Studies*, 2, 322–6.
- Dennett, D. C. (2003) Explaining the “magic” of consciousness. *Journal of Cultural and Evolutionary Psychology*, 1, 7-19.
- Eccles, J. C. (1994) *How the Self Controls its Brain*. New York: Springer.
- Elitzur, A. C. (1989) Consciousness and the incompleteness of the physical explanation of behavior. *The Journal of Mind and Behavior*, 10, 1-19.
- Elitzur, A. C. (1991) Neither idealism nor physicalism: A reply to Snyder. *The Journal of Mind and Behavior*, 12, 303-307

- Elitzur, A. C. (1994) Let there be life: Thermodynamic reflections on biogenesis and evolution. *Journal of Theoretical Biology*, 168, 429–459.
- Elitzur, A. C. (1996) Consciousness can no more be ignored: Reflections on Moody's Dialogue with Zombies. *Journal of Consciousness Studies*, 2, 353-358.
- Elitzur, A.C., Dolev, S., & Kolenda, N., Editors) *Quo Vadis Quantum Mechanics?* New York: Springer.
- Fanchon, E., Neori, K-H., and Elitzur, A. C. (2009) What does Maxwell's demon select, and how? Preprint.
- Flanagan, O. (1997) Conscious inessentialism and the epiphenomenalist suspicion. In Block, N., et al. [Eds.] *The Nature of Consciousness: Philosophical Debates*, pp. 357-373.
- Flanagan, O., and Polger, T. W. (1995) Zombies and the function of consciousness. *Journal of Consciousness Studies*, 2, 313-321.
- Griffin, D. R. (1992) *Animal Minds*. Chicago: University of Chicago Press.
- Hodges, A. (1988) Alan Turing and the Turing machine. In Herken, R. [Ed.] *The Universal Turing Machine, a Half-Century Survey*, Oxford: Oxford University Press.
- Jackson, F. (1986) What Mary didn't know, *Journal of Philosophy* 83, 291-295.
- Jahn, R. G., and Dunne, B. J., (1987) *Margins of Reality: The Role of Consciousness in the Physical World*. New York: Harcourt Brace Jovanovich.
- Jahn, R. G., and Dunne, B. J. (2001) A modular model of mind/matter manifestations (M5). *Journal of Scientific Exploration*, 15, 299-329.
- Kim, J. (1998) *Mind in a Physical World: An Essay on the Mind-Body Problem and Mental Causation*. Cambridge (Mass.): MIT Press,.
- Ludlow, P., Nagasawa, Y., and Stoljar, D., Eds. (2004) *There's Something About Mary*. Boston: MIT Press.
- Moody, T. (1994) Conversations with zombies. *Journal of Consciousness Studies*, 1, 196–200.
- Nagel, T. (1974) “What is it like to be a bat?” *Philosophical Review*, 83, 435-450.
- Penrose, R. (1994) *Shadows of the Mind: A Search for the Missing Science of Consciousness*. Oxford: Oxford University Press.
- Ridgway, P. (1996) Sleepwalking: Insanity or automatism?, *E-Law: Murdoch University Electronic Journal of Law*, 3, No.1.
- Stapp, H. P. (2009) *Mind, Matter and Quantum Mechanics*. 3rd Edition. New York: Springer.
- Tuszynski, J., Ed. (2006) *The Emerging Physics of Consciousness*. New York: Springer.
- Wilkes, K. (1984). Is consciousness important? *British Journal for Philosophy of Science*, 35, 223-243