

## EVOLUTIONARY EXPLANATION AND THE HARD PROBLEM OF CONSCIOUSNESS

Steven Horst, *Department of Philosophy, Wesleyan University,  
Middletown CT 06459, USA. Email: shorst@wesleyan.edu*

---

**Abstract:** Chalmers and others have argued that physicalist microexplanation is incapable of solving the 'hard problem' of consciousness. This article examines whether evolutionary accounts of the mind, such as those developed by Millikan, Dretske and Flanagan, can add anything to make up for the possible short falls of more reductionist accounts. I argue that they cannot, because evolutionary accounts explain by appeal to a selectional history that only comes into the picture if consciousness can first arise due to spontaneous mutation in some individual organism, and explaining this emergence of consciousness from DNA and embryology calls for precisely the kind of structurally-based supervenience account that Chalmers *et al.* have objected to. Not only does an evolutionary account not succeed where a reductionist account fails; the evolutionary account presupposes the possibility of a reductionist account.

---

David Chalmers (1995; 1996) has recently done philosophy the favour of distinguishing the 'hard problem of consciousness' — why it is that conscious phenomena appear in the world at all — from the 'easy' problems such as the ability to discriminate, categorize, and react to environmental stimuli and the focus of attention. (One assumes that the choice of the word 'easy' for these problems is intended to be somewhat droll, and true only by comparison with the hard problem in which he is primarily interested.) Chalmers argues that the hard problem cannot be solved in physicalist terms, and suggests that consciousness be viewed as being or involving a distinct kind of fundamental property in addition to those required for basic physics. Chalmers' arguments are directed against all attempts to explain the phenomenological, experiential, first-person side of consciousness in physical terms — those that do so by way of neuroscience as well as those that would try to do so directly from physics; those that appeal only to properties internal to the conscious being and those that appeal to relational (physical) properties.

I happen to think that his arguments are successful, as are those of Kripke (1971), Nagel (1974), Searle (1992), Jackson (1982), and Horst (1996). However, that is not the subject of my paper. Instead, I wish to examine whether a different form of naturalistic explanation — in this case, explanation in teleofunctional, evolutionary terms — can succeed where explanation in physical terms is seen to fail. If one agrees that we cannot answer the hard problem of consciousness in physical terms, this can be viewed as an examination of whether evolutionary explanation can save naturalism about the mind. If one is not yet convinced that consciousness cannot be explained in physicalist terms, what follows may be viewed more modestly as an examination of whether evolutionary explanation of consciousness can contribute anything to the solution of the hard problem not already contained in more structurally-based forms of physical explanation.

What I shall argue is that evolutionary explanation does *not* provide a solution to the hard problem: indeed, it would need to be supplemented by a more traditional physicalist account to do so, and hence contributes nothing towards the solution of the hard problem, in spite of being a viable and useful form of explanation with more modest virtues.

### Teleofunctions and Evolutionary Explanation

Would-be naturalizers of the mind have taken as their models a number of different paradigms from the natural sciences. Perhaps the most influential of these traces its roots to the Galilean method of resolution and composition, according to which explaining a phenomenon involves breaking it down into its component parts and then demonstrating how the behaviour of the parts necessarily produces the behaviour of the whole. Classical reductionism, type physicalism and local supervenience accounts are all inspired by this model, which in turn is styled upon the paradigm of geometric proof and construction. Recently, however, there have been two important kinds of move away from the Galilean model in philosophy of mind. The first is the growing movement towards *externalism*, in which things going on outside of the organism can play a role in determining the nature of its mental states. The second is the re-emergence of approaches to the mind drawing upon the paradigm of the Darwinian revolution in biology in the work of writers like Millikan (1984), Papineau (1993), Dretske (1995) and Flanagan (1992), among others. What is distinctive about this form of explanation is that appeals to the *function* of a phenotypic feature of an organism play a crucial role in the explanation of that feature, and the presence of that functionally-characterized feature is explained historically through a process of natural selection.

I should stress here that the notion of 'function' that is involved in such explanations is, in a broad sense, a teleological notion — in very rough, pre-theoretical terms, the function of a phenotypic feature is the selective advantage conferred upon the organism or upon the population bearing the gene for that feature. This use of the word 'function' should be carefully distinguished from the mathematical notion that is used in machine functionalist views of the mind.

There are differences in the details of how writers who champion evolutionary approaches to the mind try to explain mental features, but I think that these are by and large irrelevant to the line of investigation to be pursued here. The problem lies not in the details of specific accounts, but in the general lineaments of evolutionary explanation, and what evolutionary accounts are and are not suited to explaining. I shall therefore give a schematic account of how evolutionary explanation proceeds, first in biology and then in psychology.

Evolutionary explanation involves two mechanisms: variation (or mutation) and selection. A phenotypic feature first appears in a population through a process of mutation, which is generally understood by contemporary biology to be a random process. Most mutations are harmful, many are fatal. Some, however, confer advantages for their possessors in the biologically crucial task of passing on one's genes. This may consist either in advantages conferred upon the individual organism that increase its chances of surviving long enough to breed (by increasing the chances of longevity, by making it more likely that it will breed sooner, or by increasing the statistical success of mating producing viable offspring) or in advantages that increase the chances of survival of the gene in the offspring (by increasing the number of the offspring or increasing the chances of their viability, say by making them less attractive to predators, shortening the gestation period or increasing parental vigilance). The process of selection is one in which statistical forces operate to increase

the chances that more adaptive phenotypic traits will endure and proliferate. As a result, the explanation of a phenotypic trait T will be of the following form:

T is present in population P because

- (a) T was produced by way of spontaneous mutation in one of the ancestors of P, and
- (b) T conferred upon its possessors in the ancestors of P selective advantage A (e.g. it enabled them to run faster, or detect bugs under tree bark undetectable to their cousins who lack T, or produce more offspring).

It should be noted that not all phenotypic features are products of selection. Some are free riders carried on the same gene with traits that were selected for. Others may have become widespread for no reason connected with selective advantage. It is simply that evolutionary biology does not provide the right conceptual machinery to account for them.

In the case of psychology, then, evolutionary explanation will treat kinds of mental phenomena — whether faculties or kinds of mental state and process — as phenotypic traits of the organisms that possess them, and will attempt to explain them in terms of the selective advantage that accounts for their proliferation. In many cases, this will be a normal biological explanation, spanning over generations within a species, and associating the phenotype with a specific genetic basis. But this form of explanation can also be adapted to shorter-term adaptations within an individual organism, as learning and perception, for example, can be viewed as more rapid processes involving spontaneous variation, adaptation and selection (cf. Sayre, 1986; Millikan, 1984). The *function* of a psychological phenomenon is then understood in terms of the selection history. To say that a certain cell in the frog's visual system is a 'bug-detector' is not so much to say something about what that cell does in this particular frog (perhaps its bug-detector is damaged, or it is never exposed to flying insects), as to report the operation performed by cognate cells in its progenitors which made them more viable frogs: in this case, allowing them to efficiently detect flying insects in their visual field and eat them. Likewise, one might view the function of pain experiences in terms of the detection of tissue damage or immediate threats to bodily integrity, and perception in terms of the detection of (salient) objective features of an organism's environment.

### Existence of a Phenotype and Teleofunctional Essentialism

Before passing on to the topic of consciousness, I should pause to note that there are two different things about a phenotype that one might wish to explain in evolutionary terms. The first is its *existence* in a given organism or population: *why* do woodpeckers have long pointed bills? Because their ancestors who had the gene for long pointed bills were better able to feed themselves, hence survive and reproduce, than those cousins with shorter or duller bills. *Why* do animals have kidneys? Because those of their ancestors that developed renal systems were better able to eliminate harmful wastes within their bodies, hence better able to survive and breed than cousins that could not do so because they were being poisoned by their own waste. The other thing that one might wish to explain, however, is the *nature* of a phenotype, by way of a form of biological or teleofunctional essentialism: *what is* that thing on the

front of the woodpecker's face? It is a tool for extracting insects from beneath tree bark. What are those things hooked up to the bladder? They are devices for extracting impurities — *because that is the action whose adaptive advantage accounts for their proliferation*. (And hence your kidneys are still kidneys even if they do not in fact serve this purpose due to some form of kidney disease — what they are is determined, not by what they in fact do, but by what the phenotype of which they are tokens was selected for.) I wish to handle these topics separately in this paper. The question that is most obviously relevant to the hard problem of consciousness is that of whether evolutionary explanation can account for the existence of consciousness; thus I shall address that topic first. I shall then examine whether evolutionary accounts of the *nature* of consciousness can make up for any shortfall in accounts of its existence that do not appeal to teleofunctional essentialism.

### Dretske's Biological Explanation of Consciousness

Let us now consider an example of biological explanation in psychology. I shall use Fred Dretske's work as an example, largely because it is simpler and hence more easily presented than other biological theories such as that of Ruth Millikan. I do not believe that anything essential to my examination trades upon the differences between accounts, unless it is perhaps the fact that Millikan (1984) eschews any claims to explain consciousness, while Dretske (1995) offers a representational account of both cognition and consciousness, and claims that the two notions turn out to be closely linked.

According to Dretske, to have a thought about a thing is to have a mental representation of it. Representation, in turn, is cashed out in terms of two notions: *indication* and *function*. A indicates B if A carries the information that B is present. But not all cases of indication are cases of representation. A represents B only if A has the function of indicating B. In the case of artefacts like writing and speech, the function is *conventional* in origin, and depends upon the actions of agents. But in the case of organisms, the origin of the function is *natural*, and is cashed out in biological terms, as outlined above. Both sensations and thoughts, for Dretske, are cases of natural representation, the difference between them deriving from a distinction between 'systemic' and 'acquired' indicator functions, respectively (Dretske, 1995, pp. 6–19). These are denoted representations<sub>s</sub> and representations<sub>a</sub>.

Consciousness, for Dretske, turns out to be closely related to natural representation. A conscious state is simply a state through which we are conscious or aware of something — Dretske uses the terms 'conscious' and 'aware' as synonyms (p. 98) — and 'seeing, hearing, smelling, tasting, and feeling are specific forms — perceptual forms — of consciousness; consciousness is the genus; seeing, hearing, and smelling are species' (p. 99). In short, 'conscious states are natural representations — representations<sub>s</sub> in the case of experiences and representations<sub>a</sub> in the case of thought. Conscious creatures are creatures in whom such states occur' (p. 104). Here we have an account of the nature of both *state consciousness* (i.e. the sense in which a mental state is said to be a conscious state) and of *creature consciousness* (the sense in which a being is said to be conscious) that depends upon a teleofunctional notion of representation.

This characterization ‘yields a plausible and natural answer to questions about the function and purpose of consciousness’ (p. 116). And this is an important question for Dretske:

If some mental states and processes are conscious, others not, one can ask . . . whether conscious ones are more effective than unconscious ones. What is the point, the biological advantage, of having conscious states and processes? Those that are conscious must differ in some relevant way from those that are not. If this is not the case, then, as Davies and Humphrey (1993, pp. 4–5) conclude, too bad for consciousness: ‘Psychological theory need not be concerned with this topic’ (Dretske, 1995, pp. 116–17).

The answer, given Dretske’s characterization of consciousness, is fairly straightforward. Animals need perception to do such things as find mates and food and avoid predators, and on Dretske’s theory, consciousness goes hand in hand with perception: ‘Take away perception — as you do, according to the present theory, when you take away conscious states — and you are left with a vegetable’ (p. 118). However, blindsighters and people with various kinds of agnosias can enjoy informational sensitivity while lacking the *experience* normally associated with perception. If the same results could be achieved without experience, why is perception accompanied by experience? Dretske’s answer (admittedly only a sketch of a much fuller answer that would need to be supplied by detailed scientific research) is that persons and animals with these kinds of deficits do not in fact have all of the same abilities to negotiate their environments as do conspecifics without the deficits, and hence ‘it remains clear that people afflicted with these syndromes are always “deeply disabled”’ (p. 121). And thus

there seems to be no real empirical problem about the function, or at least a function, of sense experience. The function of sense experience, the reason animals are conscious of objects and their properties is to enable them to do all those things that those who do not have it cannot do. This is a great deal indeed. If we assume . . . that there are many things people with experience can do that people without experience cannot do, then *that* is a perfectly good answer to questions about what the function of experience is. That is why we, and a great many other animals, are conscious of things. Maybe something else besides experience would enable us to do the same things, but this would not show that experience didn’t have a function. All it would show is that there was more than one way to skin a cat — more than one way to get the job done. It would not show that the mechanism that did the job didn’t have the function of doing it (Dretske, 1995, pp. 121–2).

I include this extended quote because of the way it nicely spells out the extent of Dretske’s commitment to biological explanation, and how it is supposed to work.

### **Does Biological Explanation Explain Consciousness?**

Now I wish to address two questions about this account of Dretske’s. First, *is it any kind of explanation at all?* And second, *if so, does it solve the hard problem of consciousness?* The first question rears its head because there is a long tradition (among proponents of mechanistic explanation) of casting doubt upon explanations that turn upon teleological notions like function. However, it should be apparent that biological explanation is good at explaining some things, even if it does not explain the same things that mechanistic explanation explains. Biological explanation can explain why a phenotypic trait is present in an individual or a species, provided that there are viable stories to be had about (a) the emergence of that trait in at least one individual

through a process of variation, and (b) the survival and proliferation of that trait in a population through the conferral of selective advantage (adaptedness) upon its bearers relative to other members of the population. In point of fact, biological explanation seldom actually produces an account of the process of mutation that leads to the initial appearance of the trait. This is so for two reasons: first, these processes are believed to be random, hence anomic, and hence not subject to special explanations. (A Lamarckian theory, by contrast, would require more in the way of explanations of mutation, as it regards these as non-random.) Second, the kinds of explanation that would be needed here — a biochemistry and/or biophysics of DNA-change and particularly the embryological explanation of how particular DNA sequences produce particular phenotypic features — are largely beyond the scope of current science. However, for purposes of explanation of species change, it is generally regarded as a harmless idealization to leave the mechanisms underlying mutation and embryology unspecified.

The exceptions, of course, are cases where there is reason to regard the production of a particular phenotype, or a particular change in phenotype, by these methods as problematic. One would be suspicious, to say the least, of a biological explanation that depended on the idea that any mutation could produce within an animal an organ that served as a perpetual motion machine, because one has reason to doubt that there can be perpetual motion machines. And likewise catastrophist theories of evolution, which countenance the possibility of mutations from, say, reptiles to birds or mammals in a single mutation, have come under suspicion because it is hard to see how there could be a mechanism that would produce such changes all at once, or do so in identical ways in a sufficient number of offspring to sustain a breeding population.

Thus, what evolutionary explanation really explains is the proliferation of a phenotype, *given the plausibility of its initial appearance*. The initial appearance is treated as something that can plausibly be attributed to random processes of mutation, ultimately to be explained by biochemistry or biophysics and embryology. *Given these assumptions*, selection tells a useful and genuinely explanatory story about the function and proliferation of the phenotype — and arguably a story that cannot be told in mechanistic terms.

What, however, does this contribute towards the solution of the hard problem of consciousness? The answer, I think, is *very little*. For what such a theory can give an account of is why consciousness would *flourish* — *given that it has appeared in the first place*. And this seems quite reasonable — creatures that are conscious are likely to have great adaptive advantages over those that are not, and particular forms of consciousness are likely to confer particular kinds of adaptive advantage. All of that seems correct insofar as it goes. But what this does not do is explain how *consciousness comes upon the scene at all*: it does not tell us (a) how the mutation that first conferred consciousness came about, or (b) how some feature of DNA gives rise to consciousness in beings who possess it. (Even if this explanation is divided into an explanation of how DNA gives rise to physiological structure and physiological structure to the capacity for consciousness, the problem is not lessened.) In short, there is nothing about the specifically evolutionary or *selectional* side of the story that sheds any light upon the existence of consciousness — about how particular biological properties might be the right sort of thing to produce consciousness in the first place — and this is precisely where the hard problems lie.

Moreover, note that kinds of explanations that would need to be supplied by biochemistry, biophysics and embryology are physical and structural explanations, and precisely the kind of explanations that Chalmers *et al.* have called into question with respect to the hard question. As a result, this is *not* one of the cases in which it is safe to treat the emergence of phenotype through spontaneous variation as a harmless idealization: there is reason to doubt that physical properties determine consciousness, and therefore there is reason to doubt that the mechanisms underlying biological mutation could produce consciousness in the first place. Selectional explanation could explain the presence of consciousness in us *given* the assumption that it appeared in our ancestors through random mutation. But if no DNA structure could determine (the capacity for) conscious experience, then the selective story never gets off the ground.

Let me draw an analogy to make the point absolutely clear. Suppose someone conjectured that some species of animal was powered by a perpetual motion machine. One can certainly see how such a power source would be to an animal's advantage: it would not need to take in energy through nutrition to replenish itself, and hence would not be subject to certain hardships of privation that would imperil those around it. In short, this is a phenotypic feature that would be highly likely to proliferate, and a selective story would be easy to tell about it. The problem is that no biological mutation can produce a perpetual motion machine, and hence the selective advantage it *would* confer (however large) cannot be explained in this way. Likewise, if (as Chalmers *et al.* have argued) no structural, or more broadly, *physical* properties of organisms can determine consciousness, the evolutionary story does not get off the ground, since the selective story *never* has *anything* to say about how a phenotype *first* appears, or how it is derived from genotype. (Of course, one might have selectional stories to tell about non-physical traits that could be passed on as well, but this would not be a form of strictly biological — or naturalistic — explanation.)

The moral of the story should be clear: an evolutionary story about consciousness can explain consciousness only if there is a story about mutation and embryology which shows how the physical properties of a genotype can give rise to the phenotype in at least one individual. The selectional story contributes nothing to *this* explanation, but only explains the survival and proliferation of phenotypic features that have already appeared on the scene. In short, *a naturalistic evolutionary story about consciousness presupposes a physicalist story about the emergence of the phenotype somewhere in the history of the species.* If physicalist theories cannot address the hard problem, evolutionary theories will provide the naturalist no solace.

### Two Objections

However, one might intervene here with two objections. First, this critique has only considered evolutionary explanation of the existence of consciousness in isolation from evolutionary and teleofunctional accounts of its *nature*. Perhaps these might allow us to circumvent the problems developed above. Second, there is an important disanalogy between the explanation of consciousness and the explanation of a purported perpetual motion organ — namely, that we know that the former exists at least as surely as we know that the latter does not exist. As a result, perhaps we are

entitled to treat the assumption of the existence of consciousness as a harmless abstraction after all. I shall address these objections in order.

The first objection might go as follows: there are differences between two types of evolutionary explanation. One type of explanation explains features that are not themselves defined in teleofunctional terms by selectional history. The second type explains features that are themselves teleofunctionally-defined by reference to their selectional history. The second type of explanation looks in some ways like a definition or a tautology, as the very same features that make a feature an F explain the proliferation of Fs. Now if one is a teleofunctional essentialist about psychological kinds, the very nature of consciousness is to be understood as that of a feature whose essential properties consist in the function it was selected to perform — e.g. enabling the animal to see, hear, smell, etc. objects in its environment. In short, consciousness simply consists in whatever faculty it is that confers these abilities, and its essential property is that of conferring them.

Now how does this affect the hard problem of consciousness? It really depends on what is included in the biological function of consciousness. Is it part of the function of consciousness to do the things that it does *in a way that involves the phenomenological properties* that are the subject-matter of the hard problem, or does it treat these as non-essential concomitants? Let us consider the first case first: the function of consciousness is understood in terms both of what it allows the organism to do (see, hear, etc.) and how it does it (namely, in a way involving a phenomenology). On this use of the word 'consciousness', it picks out a feature that essentially has a phenomenology, even though there might be other mechanisms that give the same informational sensitivity without it. (In Ned Block's terminology (1995), it incorporates both *access* consciousness and *phenomenal* consciousness.) If this is what the teleological essentialist means, the criticisms I advanced earlier are untouched — for the problem remains of how mutations in DNA could produce a state having these properties.

If, on the other hand, the definition of 'consciousness' is narrowed to include only the conferral of adaptive advantage and excludes the phenomenological concomitants, the problem shapes up differently. In this case, the selectional story does not call for a prior physicalist story that explains the emergence of the phenomenology, because the phenomenology does not enter the selectional and teleofunctional story. This, however, does not so much solve the hard problem as ignore it. In this case, the teleofunctional essentialist is simply using the word 'consciousness' in a different way from the way it is used by those interested primarily in the phenomenology. - Better to say that the evolutionary theorist has explained 'focal cognition' or 'attention' but ignored phenomenology. But this does not mean that phenomenological features are not real, even if they are not subject to biological explanation. (Compare: not all phenotypic features of animals are products of selection, but they are no less real as a result.) The teleofunctionalist may give up on the project of giving a teleofunctional account of phenomenology — and may even be right to do so — but the problem of explaining such features, and hence the hard problem of consciousness, does not go away as a result. It may well be that phenomenal consciousness is not a property that confers selective advantage on its bearers. Or, alternatively, it may be that it does confer this advantage but its presence cannot be accounted for by biochemistry, biophysics and embryology. But even if Davies and Humphrey are right that this means that 'psychological theory need not be concerned with this topic' (1993, pp. 4–5), this

does not mean that the hard problem goes away. It merely means that it is not solved by psychological theory and its solution is not needed for scientific psychology to proceed apace (cf. Horst, 1996, chapter 11).

Finally, let us return to the analogy drawn earlier between biological explanations of consciousness and of perpetual motion organs. The analogy consisted in the fact that, however straightforward the *selectional* story for such features would be (because both traits would confer clear advantages upon their possessors), the biological explanation would be imperiled by the implausibility of explaining the appearance of the trait in the first place in physical terms (in the form of biochemical, biophysical and embryological stories). But of course there is an important disanalogy between the two cases as well: we know (or at least have very strong reason to believe) that no physical device can be a perpetual motion machine, and we know at least equally well that there are such things as conscious states in human beings, particularly in one's own case. There is a difference between assuming that a thing that definitely does exist is a product of a mutational process and assuming the same thing about a thing that definitely does *not* exist!

This is, of course, correct. But it does not damage my argument in the slightest. No one who thinks that there is a hard problem of consciousness believes that consciousness does not exist. Nor do they dispute that, *if* one could explain consciousness in physical terms, one would thereby have an explanation of it that could be further exploited in an evolutionary theory. The problem lies in the combination of facts that: (1) there seems to be a problem with providing a physicalist explanation of phenomenology, (2) the selectional side of the biological story presupposes the explainability of the emergence of a trait in biochemical, biophysical and embryological terms, and (3) this explanation would necessarily be a physicalist explanation. In short, no biological explanation of the hard problem of consciousness is more plausible than the physicalist explanation of the initial appearance of a trait, because the selectional story contributes exactly nothing to the solution of this particular problem, but only tells why such traits would proliferate once they appeared. And the analogy between the lack of biological explanations of consciousness and of perpetual motion machines is a fairly robust one if one focuses it in the right way. In both cases we are asking whether a *particular type of explanation* is of the right sort to account for a possible trait; and in both cases, the answer would appear to be *no*.

### Conclusion

I should repeat that I have not attempted here to argue that the hard problem of consciousness cannot be given a physicalist solution. What I have argued, rather, is that, *if* it cannot be given a physicalist solution, it cannot be given a solution in teleofunctional biological terms either. The reader who is persuaded by the arguments of writers like Kripke, Nagel, Searle, Jackson, Chalmers and Horst may thus read this as a refutation of the possibility of naturalizing consciousness in Darwinian terms. The reader who is not thus persuaded may see it as a reduction of questions about two kinds of naturalism to a question about a single kind.

