Chapter 3

of Brain States Reduction, Qualia, and the Direct Introspection

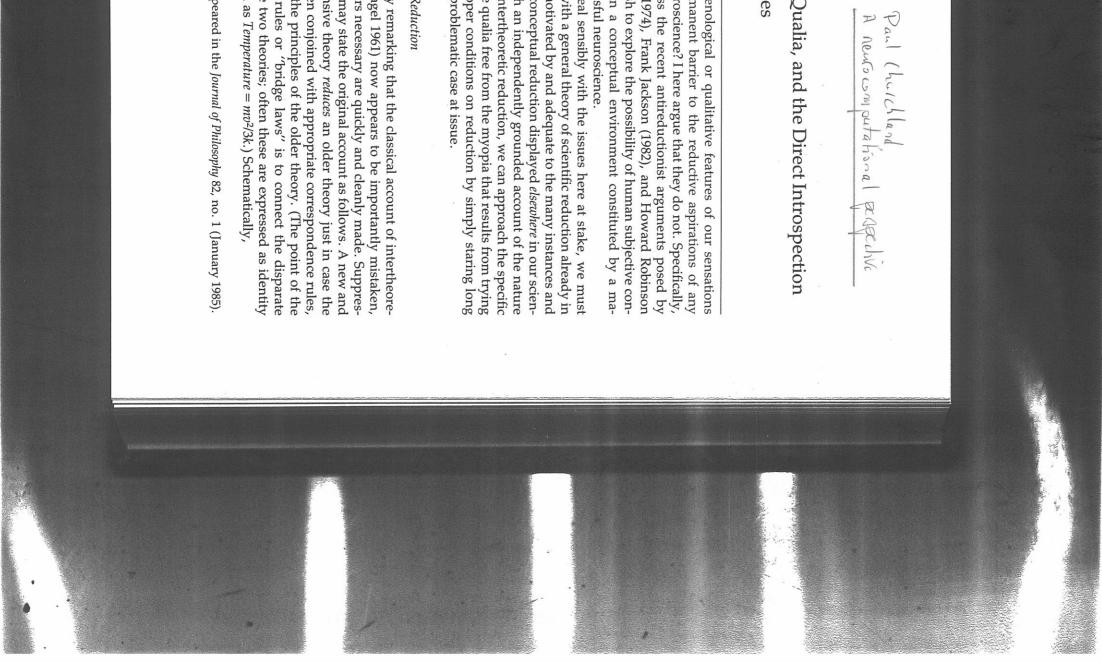
I wish to address the recent antireductionist arguments posed by Thomas Nagel (1974), Frank Jackson (1982), and Howard Robinson (1982). And I wish to explore the possibility of human subjective contured and successful neuroscience. sciousness within a conceptual environment constituted by a mamaterialistic neuroscience? I here argue that they do not. Specifically, constitute a permanent barrier to the reductive aspirations of any Do the phenomenological or qualitative features of our sensations

approach them with a general theory of scientific reduction already in and grounds of intertheoretic reduction, we can approach the specific tific history. With an independently grounded account of the nature and hard at the problematic case at issue. to divine the proper conditions on reduction by simply staring long case of subjective qualia free from the myopia that results from trying varieties of interconceptual reduction displayed elsewhere in our scienhand, a theory motivated by and adequate to the many instances and If we are to deal sensibly with the issues here at stake, we must

Intertheoretic Reduction

statements, such as $Temperature = mv^2/3k$.) Schematically, sing nicities, we may state the original account as follows. A new and We may begin by remarking that the classical account of intertheoretic reduction (Nagel 1961) now appears to be importantly mistaken, ontologies of the two theories; often these are expressed as identity correspondence rules or "bridge laws" is to connect the disparate though the repairs necessary are quickly and cleanly made. Suppreslogically entails the principles of the older theory. (The point of the new theory, when conjoined with appropriate correspondence rules, more comprehensive theory reduces an older theory just in case the

This essay first appeared in the Journal of Philosophy 82, no. 1 (January 1985).



10.

dynamics in the second and third) be somehow false as well, in contradiction to their assumed truth. theories (statistical thermodynamics in the first case, Newtonian us tollens would thus require that the premises of the new reducing uniform, as in Galilean dynamics; etc.) If reduction is deduction, mod-Keplerian astronomy; the acceleration of falling bodies isn't really thermodynamics; the planets don't really move in ellipses, as in respects, false. (Real gases don't really obey $PV = \mu RT$, as in classical duced theories turn out to be, strictly speaking and in a variety of Difficulties with this view begin with the observation that most re-

assumptions. ceded, since it is safely confined to those limiting or counterfactual different from zero). Falsity in the reducing premises can thus be conmechanical energy, or that the mass of the planets is negligible comboundary conditions (such as that the molecules of a gas have only the premises of a reduction must often include, not just the new repared to the sun's, ducing theory, but also some limiting assumptions or counterfactual This complaint can be temporarily deflected by pointing out that or that the distance any body falls is negligibly

traditional account at issue. ontology can enjoy reduction, and this fact is problematic on the statistical thermodynamics. In sum, even theories with a nonexistent of "caloric") still finds an overtly fluid thermodynamics (i.e., one committed to the existence with, nor even coextensive with, mean molecular kinetic energy. But tion. For a second example, neither is caloric-fluid-pressure identical nian by Einsteinian mechanics is a paradigm of a successful reducwith it, even at low velocities. Nevertheless, the reduction of Newtomass is features are illusory and uninstantiated. For example, tical with, nor even nomically connected with, old features, if the old play a problematic status. Newly conceived features cannot be idensome cases the reduced theory is so radically false that some or all of This defense will not deal with all cases of falsity, however, since in ontology must be rejected entirely, and the "correspondence connecting that ontology to the newer ontology therefore disnot identical with Newtonian mass, nor even coextensive a moderately impressive reduction within relativistic

duced in a reduction is the theory to be reduced. A more accurate, Cases like these invite us to give up the idea that what gets de-

general, and illuminating schema for intertheoretic reduction is as

 $T_{
m N}$ & limiting assump. & boundary cond.

logically entails

 $I_{\rm N}$ (a set of theorems of (restricted) $T_{\rm N}$), e.g., $(x)(Ax \supset Bx)$, $(x)((Bx \& Cx) \supset Dx)$,

which is relevantly isomorphic with

 T_{O} (the older theory), e.g., $(x)(Jx \supset Kx)$, $(x)((Kx \& Lx) \supset Mx)$.

allows that a true theory might reduce even a substantially false one. not necessarily as material-mode statements, but as mere ordered play no part whatever in the deduction. They show up only later, and expressed in the vocabulary proper to T_N. The correspondence rules To itself, but rather of a roughly equipotent *image* of To, an image still That is to say, a reduction consists in the deduction, within T_N , not of the older predicate it encompasses has no extension whatever. correspondence rule is entirely consistent with the assumption that the target of a relevantly adequate *mimicry*. Construed in this way, a indicate which term substitutions in the image I_N will yield the principairs: $\langle Ax, Jx \rangle$, $\langle Bx, Kx \rangle$, $\langle Cx, Lx \rangle$, $\langle Dx, Mx \rangle$. Their function is to ples of $T_{\rm O}$. The older theory, accordingly, is never deduced; it is just

reduction is drawn from P. M. Churchland 1979, section 11. For a more detailed account, see Hooker 1981.) nificant explanatory or predictive loss. (This sketch of intertheoretic intra-theoretic deduction (of I_N within T_N) and the intertheoretic mapexplanatory and predictive resources of the reduced theory. tive resources which parallel, to a relevant degree of exactness, the new or more comprehensive theory contains explanatory and predicthe older theory can be displaced wholesale by the new without sigping (of $T_{\rm O}$ into $I_{
m N}$) jointly constitute a fell-swoop demonstration that The point of a reduction, according to this view, is to show that the

respondence rule does not itself make such a claim. At best, it records and that electric current = net motion of charged particles. But a corwaves between 0.35 μ m and 0.75 μ m, that sound = atmospheric comthe fact that the new predicate applies in all those cases where its pression waves, that temperature = mean molecular kinetic energy, Material-mode statements of identity can occasionally be made, of We do wish to assert that visible light = electromagnetic

reduction for specific properties, as opposed to entire theories, and it theory T_N just in case or conceptual framework T_{O} , is reduced to a property G in some new "emergent" properties. A property F, postulated by an older theory allows us to frame a useful conception of the contrary notion of The preceding framework allows us to frame a useful conception of

- (1) $T_{\rm N}$ reduces $T_{\rm O}$,
- (2) 'F' and 'G' are correspondence-rule paired in the reduction, and
- (3) the reduction is sufficiently smooth to sustain the ontology of T_O and thus to sustain the identity claim '*F*-ness = *G*-ness'.

the laws of T_N). laws of $T_{
m O}$) are a subset of the causal powers of G-ness (as outlined in to G-ness just in case the causal powers of F-ness (as outlined in the Intuitively, and in the material mode, this means that F-ness reduces

to T_N) just in case Finally, a property F will be said to be an emergent property (relative

- (1) F is definitely real and instantiated,
- (2) F is cooccurrent with some feature or complex circumstance recognized in T_{N_r} but
- (3) F cannot be *reduced* to any property postulated by or definable within T_N .

ments substantive questions of emergence and irreducibility in a few moant theories. Having outlined these notions, we shall turn to address fore claims about the relative poverty in the resources of certain aspirby F-ness. Claims about the emergence of certain properties are thereadequate to define a property with all of the causal powers possessed Intuitively, this will happen when T_N does not have the resources

encounters a different sense of 'emergent', one that implies that an exactly when the elements of some substrate are suitably organized, a senses. In scientific contexts, one frequently hears it used to apply to that is at issue in this paper.) the second implies irreducibility. It is emergence in the second sense feature of its substrate. The first sense positively implies reducibility; emergent property does not consist in any collective or organizational numbered among them. But in philosophical contexts one more often properties, and quite probably the qualia of our sensations should be this innocent sense of 'emergent', there are a great many emergent tain the set of causal powers ascribed to the "emergent" property. In certain relations to one another, a set of relations that collectively susproperty that consists in the elements of that substrate standing in what might be called a "network property, 'emergent property' (A word of caution is perhaps in order here, since the expression is often used in two diametrically opposed " a property that appears

tion at longish wavelengths. mean kinetic energy and will collectively emit electromagnetic radiamolecular motion, one cannot deduce from it that a roaring hearth only that water will scatter electromagnetic radiation at such and such cerning H₂O, one cannot deduce from it that water will be blue, but however much one bends and squeezes the molecular theory conof warmth or blueness must be irreducibly emergent properties, since T_N . It is occasionally claimed, for example, that the objective features that the existence or appearance of F-ness cannot be deduced from property F relative to some theory T_{N} , it is not sufficient to point out emphasized. The first is that in arguing for the emergence of a given will be warm, but only that its molecules will have such and such a wavelengths. And however much one wrings from the mechanics of Before we continue, several points about reduction need to be

condition of direct deducibility. For example, formal considerations alone guarantee that, for any predicate 'F' not already in the prowhy it would be even more foolish to insist on the much stronger saw at the beginning of this section. And there are additional reasons correspondence rules) a requirement on successful reduction, as we to make even indirect deducibility (i.e., deducibility with the help of conclusion against reducibility does not follow. It is a serious mistake any conceptual framework to reduce any other, distinct conceptual prietary lexicon of the aspirant reducing theory I_N , no statements thus trivialize the notion of reduction by making it impossible for tions) will be deducible from T_N . The deducibility requirement would whatever involving 'F' These premises about nondeducibility are entirely true, but the Even temperature, (beyond tautologies and other trivial excepthat paradigm of a successfully

reduced property, would be rendered irreducible, since the term does not appear in the lexicon of statistical mechanics

mal or color concepts is evidently this same unreasonable demand. gy for us, as well as mechanics, or electromagnetic theory, or what have you. The demand that molecular theory directly entail our therwould be to insist that the new theory do predictive cultural anthropoloidiosyncratic human culture is going to conceive of that domain. That that it also be able to predict how this, that, or the other conceptually of which would be roughly isomorphic (each in its different way) with in (that is, function successfully in) that domain. But we cannot insist some part of the correct account that a utopian theory will eventually could have been roughly adequate to common experience, and many frameworks to describe the observable world, each one of which might have used any one of an infinite number of other conceptual theory of a given objective domain that it account for the phenomena currently use precisely the conceptual framework we do use. We too strong. The fact is, it is an historical accident that we humans There is a further reason why the demand for direct deducibility is Accordingly, we can legitimately ask of a putatively correct

within a much more penetrating conceptual systemsame set of objective properties. The hypothesized identity of the be for concluding that both theories have managed to latch onto the systematic nomological parallels constitute the best grounds there can both theories presume to describe the same empirical domain, these nomically parallel: they are both at least partially correct accounts of properties at issue explains why $I_{
m N}$ and $T_{
m O}$ are taxonomically and powers/roles/reatures are systematic analogues of the powers/roles/features of the set of properties postulated by the old theory. Since powers/roles/features are systematic analogues of the resources to All we can properly ask of a reducing theory is that it have the same objective properties. I_N merely frames that account conjure up a set of properties whose nomological -that of T_N .

wholly false and illusory. epistemic stagnation or the outright elimination of old frameworks as this conceit, then the only alternatives to intertheoretic reduction are cies comprise an exhaustive account of anything at all? If we put aside assert that the feeble conceptual achievements of our adolescent spethose in turn by frameworks better still, for who will be so brash as to will eventually be reduced or displaced by new and better ones, and Moreover, it is to be expected that existing conceptual frameworks

2 Theoretical Change and Perceptual Change

only things that occasionally enjoy intertheoretic reduction. Observ-Esoteric properties and arcane theoretical frameworks are not the

identical with having a certain triplet of electromagnetic reflectance of microscopically embodied energies, and so forth. efficiencies; being warm is identical with having a certain mean level with being an oscillation in air pressure at 440 hertz; being red is able properties and commonsense conceptual frameworks can also smooth reduction. Thus, being a middle-A sound is identical

of our new perceptual judgments than we made of our old ones. new conceptual framework. We can thus make better inferential use a second reason: the greater inferential or computational power of the native perceptual equipment. Such displacement is also desirable for the old framework. We can thus make more penetrating use of our in the discriminatory reach of our native perceptual systems, though since the new vocabulary observes distinctions which are in fact withcated reducing theory. It is even desirable that we begin doing this, spontaneous perceptual reports in the language of the more sophistierty identities just listed, it is quite open to us to begin framing our those objective distinctions go unmarked and unnoticed from within in all of its observational contexts as well. Given the reality of the propshould be appreciated that the reducing theory can displace the old framework old framework not just in contexts of calculation and inference. It Moreover, the relevant reducing theory is capable of replacing the

once established. A nonscientific example may help to get the initial transformations and the naturalness of the new conceptual regime point across. It is difficult to convey in words the enormity of such perceptual

tured detail, concerning which the child is both dumb and deaf. melody line. The matured musician hears an entire world of strucguishable and identifiable chords supporting an appropriately related apprehended tune is now a rationally structured sequence of distinnow a mosaic of distinguishable elements. What was before a dimly of the orchestra performing it. What was before a seamless voice is symphony, and the same person's apprehension of the same symthe gap between an untrained child's auditory apprehension of a phony forty years later, when hearing it in his capacity as conductor Consider the enormous increase in discriminatory skill that spans

abyss, scattering nearby planets, yellow dwarf stars, blue and red tions he can estimate with accuracy. Or consider the astronomer, for whom the speckled black dome of her youth has become a visible tannin, acid, carbon dioxide, and so forth, whose relative concentratwenty distinguishable elements: ethanol, glycol, fructose, sucrose, "red wine" used by most of us divides into a network of fifteen or ticed and chemically sophisticated wine taster, for whom the category Other modalities provide comparable examples. Consider the prac-

her unaided (repeat: unaided) eye. discriminable as such and locatable in three-dimensional space with giants, distant globular clusters, and even a remote galaxy or two, all

Such frameworks are characteristically a cultural heritage, pieced examined at length in P. M. Churchland 1979, sections 1 through 6.) in their absence. (The role of theory in perception and the sysness and penetration to our sensory lives that would be impossible together over many generations, and their mastery supplies a richdomain than is immediately apparent to untutored discrimination work that embodies far more wisdom about the relevant sensory tematic enhancement of perception through theoretical progress are In each of these cases, what is finally mastered is a conceptual -whether musical, chemical, or astronomical--a frame-

discriminatory mechanisms remain unchanged recognition could be very much greater than it is, though our native nature than it actually does, our introspective discrimination and native discriminatory mechanisms remain the same. Correlatively, if ternal states and activities would be much diminished, though our connecting generalizations, our introspective apprehension of our infoundly. If it embodied substantially less wisdom in its categories and ment in its own right, and it shapes our matured introspection proordinary language is a modestly sophisticated theoretical achieveconceptual framework for psychological states that is embedded in bodied in the psychological vocabulary of the language we learn. The criminations that others are already making, the discriminations emare those it is useful for us to make. Generally, those are the disoften quite slowly. And the specific discriminations we learn to make most part learned; they are acquired with practice and experience, phenomenon. The introspective discriminations we make are for the folk psychology embodied substantially more wisdom about our inner Our introspective lives are already the extensive beneficiaries of this

sider now the possibility of learning to describe, conceive, and introthat were framed, as a matter of course, in the appropriate concepts respond to that reconfigured discriminative activity with judgments from a completed neuroscience. And suppose we trained ourselves to language, but to some more penetrating taxonomy of states drawn responded not to the primitive psychological taxonomy of ordinary make a new and more detailed set of discriminations, a set that corsense folk psychology. Suppose we trained our native mechanisms to that successfully reduces, either smoothly or roughly, our commonthe conceptual framework of a matured neuroscience, a neuroscience spectively apprehend the teeming intricacies of our inner lives within This brings me to the central positive suggestion of this paper. Con-

vironment prepared largely by Sellars 1956. The idea has been explored more recently in P. M. Churchland 1979 and in chapter 1 Rorty who first identified and explored this suggestion. See Feyerabend 1963a and Rorty 1965. This occurred in a theoretical enfrom neuroscience. (I believe it was Paul K. Feyerabend and Richard

quantum leap in self-apprehension. inferential application. But that seems a small price to pay for the science in order to pull this off. And we will have to practice its noncourse have to learn the conceptual framework of a matured neurofocus of a trained musician's auditory discrimination. We will of just as Gm7 chords and Adim chords are moved into the objective moved into the objective focus of our introspective discrimination, layer of the occipital cortex, inhibitory feedback to the lateral genicuspiking frequencies in specific neural pathways, resonances in the nth approximate a revelation. Dopamine levels in the limbic system, the a fair parallel, then the enhancement in our introspective vision could astronomer (who can see the temperature of a blue giant star) provide late nucleus, chords), the enologist (who can see and taste the glycol), and the If the examples of the symphony conductor (who can hear the Am7 and countless other neurophysical nicities could be

several neurosciences. and familiar, and it centers on the growing explanatory success of the to be settled a priori. The evidence for a positive answer is substantial theories will prove able to do this is a wholly empirical question, not generalizations of folk psychology. Whether future neuroscientific set of embedding laws that faithfully mimics the taxonomy and causal need only include, or prove able to define, a taxonomy of kinds with a logical states and properties. A matured and successful neuroscience eventual reduction of mental states and properties to neurophysio-All of this suggests that there is no problem at all in conceiving the

now examine their arguments. emergence is the correct story for our mentalistic ontology. Let us however, take a quite different line. They find no fault with folk The qualia-based arguments of Nagel, Jackson, outright elimination as the eventual fate of our mentalistic ontology. tion of its familiar ontology. That line suggests substantial revision or whether it has the categorial integrity to merit the reductive preservapsychology. Their concern is with the explanatory and descriptive tory and predictive poverty of folk psychology, and they question myself (1981a). My negative arguments there center on the explana-But there is negative evidence as well; I have even urged some of it any possible neuroscience, and their line suggests that and Robinson,

3 Thomas Nagel's Arguments

three, I shall argue, are unsound. matured neuroscience. any plausible or adequate reduction within the framework of arguments in support of the view that such properties will never find neuroscience. In his classic position paper (1974), I find three distinct stitute a problem for the reductive aspirations of any materialistic ences, the properties or qualia displayed by our sensations, that con-For Thomas Nagel, it is the phenomenological features of our experi-All three arguments are beguiling, but all

The first argument

from reductions elsewhere in science, says Nagel, is that What makes the proposed reduction of mental phenomena different

it is impossible to exclude the phenomenological features of experience from a reduction, in the same way that one excludes the phenomenal features of an ordinary substance from a physical or chemical reduction of it—namely, by explaining them as effects on the minds of human observers. (1974, p. 437)

subjective point of view. But this is not what interests me about this the phenomenological features are essential to experience, and to the substances elsewhere in science exclude the phenomenal features of the The reason it is impossible to exclude them, continues Nagel, is that What interests me is the claim that reductions of various

phenomenal properties. Despite widespread ignorance of their dywith the mean level of the objects' microscopically embodied enerobjective phenomenal property of apples, properties to which everyone's perceptual mechanisms are keyed. namical and mechanical properties, out there in the objective world, are genuine gies. Pitch, an objective phenomenal property of a sound, is identical with its oscillatory frequency. These electromagnetic and micro-Warmth, an objective phenomenal property of objects, is identical tain wavelength triplet of electromagnetic reflectance efficiencies These properties are not excluded from our reductions. Redness, an of an apple, the warmth of a coffee cup, and the pitch of a sound. phenomenal features at issue are those such as the objective redness This is simply false, microphysical details, it is these objective physical and the point is extremely important. The is identical with a cer-

monsense vocabulary for observable properties, and learn to frame plete that one can already displace entirely large chunks of our com-The reductions whose existence Nagel denies are in fact so com-

Churchland 1981b, pp. 128–130 [this volume, chapter 2, pp. 30–31].) issue. (See my 1979, sections 2 through 6. See also Paul and Patricia adequate to permit the reliable discrimination of the properties at roles of the properties thus discriminated. But they are abundantly their presence from their absence. They have been doing so for milisms can easily discriminate such properties, one from another, and 80 percent. These microphysical and electromagnetic properties can be felt, heard, and seen, respectively. Our native sensory mechanhertz. And the three critical electromagnetic reflectance efficiencies The mean kinetic energy (KE) of the molecules in this room, for example, is currently about 6.2×10^{-21} joules. The oscillatory frequency of (at 0.45, 0.53, and 0.63 μ m) of this white piece of paper are all above this sound (I here whistle C one octave above middle C) is about 524 one's perceptual judgements directly in terms of the reducing theory. to reveal the microphysical details and the extended causal The "resolution" of these mechanisms is inadequate,

dodge is no longer open to us once the problematic properties are already located within the mind. mental phenomena. And as Nagel correctly points out, the relocation only confront them again later as we address the place in nature of to try to "kick the phenomenal properties inwards," since that only the contrary, they are as objective as you please, with a wide variety of objective causal properties. Moreover, it would be a mistake even postpones the problem of reckoning their place in nature. We shall have no real existence save inside the minds of human observers. On On this view, the standard perceptual properties are not "secondproperties at all, in the standard sense that implies that they

the internal cases are not different: they are parallel after all. confronted with parallel forthrightness, and can be reduced where they stand: inside the human observer. So far then, the external and phenomenal properties are so treated, then subjective qualia can be observer. As we saw, this can and has in fact been done. If objective and they should be reduced where they stand: outside the human of observers" in the first place. They should be confronted squarely, warmth, etc.) should never have been "kicked inwards to the minds immune from the sort of reductions found elsewhere in science. I draw a Nagel concludes from this that subjective qualia are unique in being very different conclusion. The objective qualia (redness,

The second argument

experiences, the qualia of sensations, are essentially accessible from A second argument urges the point that the intrinsic character of

following argument. This somewhat diffuse argument appears to be an instance of the

- (1) The qualia of my sensations are directly known by me, by introspection, as elements of my conscious self.
- (2) The properties of my brain states are *not* directly known by me, by introspection, as elements of my conscious self.
- ∴ (3) The qualia of my sensations \neq the properties of my brain states.

to the first: And perhaps there is a second argument here as well, a complement

- The properties of my brain states are known-by-thevarious-external-senses, as having such and such physical properties.
- (2) The qualia of my sensations are *not* known-by-the-various-external-senses, as having such and such physical properties.
- \therefore (3) The qualia of my sensations \neq the properties of my brain states.

The argument form here is apparently

- (1) Fa
- (2) $\sim Fb$
- $\therefore (3) \quad a \neq b.$

cases is amply illustrated in the following parallel arguments. this is a valid argument form. But in the examples at issue, F is obviously not an extensional property. The fallacy committed in both Given Leibniz's Law and the extensional nature of the property F, this is a valid argument form. But in the examples at issue, F is

- (1) Hitler is widely recognized as a mass murderer.
- (2) Adolf Schicklgruber is *not* widely recognized as a mass murderer.
- ∴ (3) Hitler ≠ Adolf Schicklgruber.

or,

- (1) Aspirin is known by John to be a pain reliever.
- (2) Acetylsalicylic acid is *not* known by John to be a pain reliever.
- \therefore (3) Aspirin \neq acetylsalicylic acid.

or, to cite an example very close to the case at issue,

- (1) Temperature is known by me, by tactile sensing, as a feature of material objects.
- (2) Mean molecular kinetic energy is *not* known by me, by tactile sensing, as a feature of material objects.
- ∴ (3) Temperature ≠ mean molecular kinetic energy.

nection with the identity theory.) Jaegwon Kim (1967) who first identified this fallacy specifically in confore hardly grounds for concluding that 'a' and 'b' cannot be core-ferential or coextensive terms! (I believe it was Richard Brandt and show a difference in truth value for two terms 'a' and 'b' of a coreferential or coextensive term for whatever holds the place of number of intensional contexts whose distinguishing feature is that is known (perceived, recognized) by me, as an F' is one of a large of my brain state'). In logician's terms, the propositional function, some specific description or other. Such apprehension is not a genuine they do not always retain the same truth value through substitution another, equally accurate, coreferential description (e.g., (e.g., 'qualia of my mental state'), and yet fail to be recognized under same subject may be successfully recognized under one description feature of the item itself, fit for divining identities, since one and the subject item's being recognized, perceived, or known as something, under ascribed in premise (1) and witheld in premise (2) consists only in the 'qualia of my sensations' and 'property of my brainstates') The problem with all of these arguments is that the "property" Accordingly, that such a context (i.e., the one at issue) should property is there-

what he wishes to defend is the following modalized version of the cumstances. In correspondence, Thomas Nagel has advised me that urged that one's brain states are more than merely not (yet) known by argument. introspection: they are not knowable by introspection under any cirent version of the argument, which we must also consider. It may be This objection is decisive, I think, but it does not apply to a differ-

(1) My mental states are knowable by me by introspection.

 \therefore (3) My mental states \neq my brain states.

Here Nagel will insist that being knowable by me by introspection is a argument is free of the intensional fallacy discussed above genuine relational property of a thing, and that this version of the

the argument contains a false premise: premise (2). At the very least, he can insist that (2) begs the question. For if mental states are indeed commits the same error instanced below. that is, are indeed knowable by introspection, and Nagel's argument their more penetrating neurophysiological descriptions. Brain states, then we can certainly learn to think of and recognize them under states under their familiar mentalistic descriptions grained nature. And if we can learn to think of and recognize those been introspecting all along, though without appreciating their fineidentical with brain states, then it is really brain states that we have And so it is. But now the reductionist is in a position to insist that -as all of us have

- l) Temperature is knowable by tactile sensing.
- (2) Mean molecular kinetic energy is not knowable by tactile sensing.
- . (3) Temperature ≠ mean molecular kinetic energy.

and more penetrating conceptual framework in their description.) are one and the same states. One would simply employ a different exactly the same as introspecting the states of one's mind, since they the states of one's brain? (What would that feel like? It would feel why is it unthinkable that one might come to know, by introspection, come to know, by feeling, the mean KE of atmospheric molecules, our native discriminatory mechanisms are keyed to. And if one can 6.2×10^{-21} joules, for whether we realize it or not, that is the property one can learn to feel that the mean KE of its molecules is about as one can learn to feel that the summer air is about 70°F, or 21°C, so therefore have a false premise. Premise (2) is clearly the stinker. Just mean molecular kinetic energy. Since the argument is valid, it must Here the conclusion is known to be false. Temperature is indeed

al fallacy. My guess is that Nagel has profited somewhat from the second premise, to distinguish it from the second premise of the very ever, meets both conditions. in the second version, the argument is valid. Neither version, howambiguity here. For in the first version, both premises are true. And first version of the argument, the version that commits the intensionmust be careful, in evaluating the plausibility of Nagel's

in the final section of this paper. For now, let us move on. The matter of introspecting one's brain states will arise once more

The third argument

keep of subjective qualia. (see Nagel 1974, pp. 438ff.) doomed to dash themselves, drawn is that the reductive aspirations of neurophysiology are nor perhaps even imagine, what it is like to be a bat. Even total knowledge of the physical details still leaves something out. The lesson and its interaction with the physical world, one could still not know, that, no matter how much one knew about the bat's neurophysiology experiences enjoyed by an alien creature such as a bat. The claim is Nagel's paper. The leading example is the (mooted) character of the The last argument here is the one most widely associated with unrealized, against the impenetrable

directly with humans, I shall confront the problem as he formulates recent paper by Frank Jackson (1982). Since Jackson's version deals I his argument is almost identical to an argument put forward in a

.

4 Jackson's Knowledge Argument

possible states. ture and activity of the brain and its visual system, of its actual and comes to know everything there is to know about the physical strucgreatest neuroscientist, all from within this room. In particular, she manages to transcend these obstacles. She becomes the world's by means of a black/white television monitor, and being brilliant, she shades of black, white, and grey. She learns about the outside world tire life in a room that is rigorously controlled to display only various Imagine a brilliant neuroscientist named Mary who has lived her en-

all mental phenomena. of-red. Therefore, complete knowledge of the physical facts of visual of seeing a ripe tomato, what it is like to see red or have a sensation-Hence, materialism cannot give an adequate reductionist account of perception and its related brain activity still leaves something out. ences were she finally to leave her room: the nature of the experience who live outside her black/white room, and about her possible experieven imagine, about the actual experiences of all the other people But there would still be something she did not know, and could not

To give a conveniently tightened version of this argument,

 Mary knows everything there is to know about brain states and their properties.

(2) It is not the case that Mary knows everything there is to know about sensations and their properties.

Therefore, by Leibniz's Law,

(3) Sensations and their properties ≠ brain states and their properties

It is tempting to insist that we here confront just another instance of think, find at least two other shortcomings in this sort of argument. rent, entirely extensional context. Let us suppose that it is. We can, I intensional Campbell 1983) insist that 'knows about' is a perfectly transpafallacy discussed earlier, but Jackson's defenders

The first shortcoming

matter of having a representation of redness in some prelinguistic or both premises, but it is not *univocal* in both premises. (David Lewis [1983] and Laurence Nemirow [1980] have both raised this same a matter of being able to make certain sensory discriminations, or sublinguistic medium of representation for sensory variables, or to be ten in neuroscience texts; whereas knowledge in (2) seems to be a mastered a set of sentences or propositions, the kind one finds writaddressed in (2). Knowledge in (1) seems to be a matter of having seems pretty clearly to be different from the kind of knowledge in both premises. But the kind of knowledge addressed in premise (1) mine.) Jackson's argument is valid only if 'knows about' is univocal objection, though their analysis of the ambiguity at issue differs from This defect is simplicity itself. 'Knows about' may be transparent in something along these lines.

do seem very different, even in advance of a settled analysis of the will be sustained so long as the type of knowledge invoked in premise ledge by acquaintance' are possible, and the charge of equivocation based. As my alternative gloss illustrates, other analyses of 'knowsense of 'knows about', but they need not be so narrowly commit-(1) is distinct from the type invoked in premise (2). Importantly, they Lewis and Nemirow plump for the "ability" analysis of the relevant and the complaint of equivocation need not be so narrowly

knowledge each has of exactly the same thing. The difference is in the may reside not in what is respectively known by each (brain states by visual cortex but has never enjoyed a sensation of red, and a person the former, qualia by the latter), but rather in the different type of who knows no neuroscience but knows well the sensation of red, In short, the difference between a person who knows all about the

the resulting argument is a clear non sequitur. son's argument with the two different expansions suggested above, one replaces the ambiguous occurrences of 'knows about' in Jackmanner of the knowing, not in the nature of the thing(s) known. If

- (a) Mary has mastered the complete set of true propositions about people's brain states.
- (b) Mary does *not* have a representation of redness in her prelinguistic medium of representation for sensory variables.

Therefore, by Leibniz's Law,

(c) The redness sensation \neq any brain state.

they do not entail (c). Premises (a) and (b) are compossible, even on a materialist view. But

exploit this variety illegitimately: both arguments equivocate tion, perhaps hundreds of them. Jackson's argument, and Nagel's, brain uses a considerable variety of modes and media of representamore modes and media of representation than the simple storage of sentences. beyond the reach of physical science. It just means that the brain uses edge" of one's sensations in a way that is independent of the scienprecludes this. The materialist can freely admit that one has "knowl-'knows about'. And this proposition is pretty obviously true: almost certainly the tific theories one has learned. This does not mean that sensations are than having mastered a set of sentences. And nothing in materialism In sum, there are pretty clearly more ways of "having knowledge"

to account for all mental phenomena. quaintance): what it is like to see red. Dualism is therefore inadequate vision. There would still be something she did not know2 (by acthing there is to know about the ectoplasmic processes underlying plasmologist" ground all mental phenomena. Let our cloistered Mary be an "ectoit 'ectoplasm'that Jackson were arguing not against materialism, but against dualism: against the view that there exists a nonmaterial substance—call form of argument were sound, it would prove far too much. Suppose This criticism is supported by the observation that, if Jackson's this time, and let her know1 (by description) every--whose hidden constitution and nomic intricacies

show nothing, one way or the other, about how mental phenomena it exploits the same equivocation. But the truth is, such arguments might be accounted for. This argument is as plausible as Jackson's, and for the same reason:

think, is the claim that Mary could not even *imagine* what the relevant experience would be like, despite her exhaustive neuroscientific objection comes to, then there is no objection worth addressing is an intolerably strong demand on reduction, and if this is all the nomenological properties. As we saw in section 1, direct deducibility neuroscience. red' will be deducible from premises restricted to the language of It is true, of course, that no sentence of the form 'x is a sensation-ofabout the subjective qualitative nature of sensations not yet enjoyed sequences to be expected from a successful neuroscientific account of found importance for understanding one of the most exciting con-There is a further shortcoming with Jackson's argument, one of proinformation. knowledge, and hence, that she must still be missing certain crucial What the defender of emergent qualia must have in mind here, knowledge of neuroscience must leave Mary hopelessly in the dark I draw your attention to the assumption that even a utopian But this is no point against the reducibility of phe-

there is to know about the physical brain and nervous system. much one might know if, as premise (1) asserts, one knew everything tor our internal states. could follow upon a wholesale revision in our conceptual framework changes in our introspective apprehension of our internal states that particular, none of these philosophers has even begun to consider the because none of these philosophers has adequately considered how premise (2) seems plausible to Jackson, Nagel, and Robinson only This claim, however, is simply false. Given the truth of premise (1),

stantial success, even in advance of receiving external stimuli that ot-grey', or 'a sensation-of-white'; but rather identifies them more conceptualize her inner life, even in introspection, perceptual domain. In particular, suppose that Mary has learned to perception through the systematic reconceptualization of the relevant sensations. would give Mary detailed information about the qualia of various would actually produce it. imagine being in the relevant cortical state, and imagine it with subconcepts for the sensational states at issue (namely, sensations-ofpital cortex (or whatever). If Mary has the relevant neuroscientific revealingly as various spiking frequencies in the nth layer of the occiher visual sensations crudely as 'a sensation-of-black', 'a sensationcompleted neuroscience we are to imagine. So she does not identify The fact is, we can indeed imagine how neuroscientific information but has never yet been in those states, she may well be able Recall our earlier discussion of the transformation of in terms of the

conceptualized clearly by her, but not previously enjoyed. succeed in this, and do so regularly on similar tests for other states, to us), and see if she can identify it correctly on introspective grounds ing frequency of 90 hertz in the gamma network: a "sensation-of-red" that would (finally) produce in her the relevant state (namely, a spik-One test of her ability in this regard would be to give her a stimulus as 'a spiking frequency of 90 hertz, the kind a tomato would It does not seem to me to be impossible that she should

skill beyond all possibility for Mary? meets the description. Skilled musicians can do this. Why is a similar which is the target, and see if he can pick it out as the sound that him brood for a bit. never have heard before and certainly does not remember. Specify for can construct, in auditory imagination, the sound of a chord he may sound in auditory imagination. Moreover, a really skilled individual him a relatively unusual onely, a trained pianist or guitarist can identify the chord and recall its chord. And the reverse is also true: if a set of notes is specified verbalabsolute pitch, one can even name the notes of an apprehended groups of discriminable notes. If one is sufficiently practised to have musical education changes this, and one comes to hear chords as able one from another, but without elements or internal structure. the young and unpracticed ear hears as undivided wholes, discriminwill show that it is not. Musical chords are auditory phenomena that This may seem to some an outlandish suggestion, but the following Then play for him three or four chords, one of -an F#9thadd13th for example -and let

chords are audibly structured sets of elements. Sensations-of-color it is tempting to reply, musicians can do this only because

dence to suggest that our sensations-of-color are indeed structured sets of especially since there has recently emerged excellent empirical evithis possibility out, and it is difficult to see how he can do this a priori, formed inspection? Jackson's argument, to be successful, must rule nal structure, unnoticed so far, but awaiting our determined and init be unthinkable that sensations-of-color possess a comparable interments. They also seemed to be undifferentiated wholes. Why should But neither did chords seem, initially, to be structured sets of ele-

tively responsive. Since colors are apprehended by us, it is a good those wavelengths to which the retina's triune cone system is selecvertices—its reflectance efficiencies at three critical wavelengths, system as being uniquely specified by its joint position along three Land (1977) represents any color apprehendable by the human visual The Retinex theory of color vision recently proposed by Edwin