# Physically orthodox theories-of-consciousness must predict biophysical mechanisms for the integration of spatially-distributed and temporally-extended neural codes

Nicholas Rosseinsky

# Physically orthodox theories-of-consciousness must predict biophysical mechanisms for the integration of spatially-distributed and temporally-extended neural codes

**Nicholas M. Rosseinsky[1,*]**

[1]Department of Neuroscience, Center for Dialog in Science, Exeter, U.K.

**\* Correspondence:** Nicholas M. Rosseinsky, Department of Neuroscience, Center for Dialog in Science, 14 Rosebarn Avenue, Exeter, EX4 6DY, U.K.

rosseinsky.nicholas.m@gmail.com

**Word count:** 11,999

## Abstract

The characterization of certain spatiotemporal brain dynamics as "encodings" *e.g.* of the external environment is central to modern neuroscience. It is usually assumed that *any* encoding scheme employed in the brain can be utilized as the informational basis for human conscious experience, thus supporting a default view that no "new physics" is required for the explanation of consciousness. Here, in contrast, it is shown that the assumption of orthodox physics severely *constrains* encoding schemes for conscious experience, specifically because physical orthodoxy requires local causality, and proposes that no physical structures beyond the brain itself are involved in the generation of consciousness. Under these conditions, spatially-distributed codes cannot be the final encoding basis for consciousness, despite their ubiquitous appearance in the current neuroscience of brain structure and function, because required aggregation of information encoded at distinct spatial locations would be definitively non-local. Similarly, rate codes and other, more complex, temporally-extended schemes are non-local in time. Thus, the assumption of physical orthodoxy limits the encoding basis for components-of-consciousness to instantaneous states at single locations. Consequently, if physically-orthodox theories are to explain consciousness, biophysical mechanisms must exist to reduce well-established spatially- and temporally-extended codes to spatiotemporally point-like states, thus creating a strong and specific prediction for the biophysical architecture of neural information processing in the human brain. Future empirical exclusion of required biophysical mechanisms would establish that physically-novel mechanisms govern the generation of consciousness. Conversely, future empirical discovery of required mechanisms would be a strong pointer to the physical basis of human consciousness, if these mechanisms have no other functional role in human biology.

## 1. Introduction

Understanding consciousness is sometimes characterized as one of the greatest outstanding scientific challenges (Miller, 2005). The mainstream approach (Anderson, 1972; Baars, 1988; Edelman, 1989; Crick and Koch, 1990; Tononi and Edelman, 1998; Tegmark, 2000; Churchland, 2005; Baars and Edelman, 2012) is that certain brain dynamics are more or less synonymous with conscious experience, thus reducing the challenge to the empirical characterization (Dehaene, 2014) of those

42  dynamical states that are both necessary and sufficient for particular conscious experiences. In this
43  mainstream view, consciousness is either identical to particular brain dynamics (Place, 1956; Smart,
44  1959; Tegmark, 2000) or a conventionally-emergent property (Anderson, 1972; Van Gulick, 2001;
45  Chalmers, 2008), so that consciousness can be explained in a physically orthodox setting, without
46  any need for "new physics" (Edelman, 1989). Given the success of modern physical theory in
47  explaining other natural phenomena within a single conceptual framework, the assumption that
48  consciousness is a physically orthodox phenomenon in this sense is perhaps a natural hypothesis.
49  Nevertheless, at present, this hypothesis has no empirical support, with the result that alternative,
50  non-orthodox, proposals [*e.g.* (Penrose, 1989; Stapp, 1993; Chalmers, 1996; Beck and Eccles, 2003)]
51  can arguably claim a scientific status equivalent to that of the mainstream view.

52  The present paper is the first in a series of three (Rosseinsky, 2014a, 2014b) whose goal is to
53  develop objective empirical tests in ordinary neuroscience that can either falsify the "physically
54  orthodox" hypothesis or provide new and significant information concerning the physical basis of
55  consciousness (in orthodox settings). Unlike the vast majority of consciousness-related work in
56  neuroscience (Tononi and Koch, 2008) that focuses on relatively macroscopic hypotheses *e.g.*
57  concerning brain areas (Merker, 2007; Shulman *et al.*, 2009) or global phase relationships (Klimesch
58  *et al.*, 2001; Melloni *et al.*, 2007), the present series focuses on *neural codes* that act as the
59  informational basis of consciousness. Requirements for codes are necessarily requirements for neural
60  wiring and associated biophysical machinery, and therefore result directly in empirically-observable
61  predictions for the class of physically orthodox theories-of-consciousness. (Because these empirical
62  tests examine the biophysical structure of the human brain, they do not depend at all on detailed
63  subjective report.)

64  This paper examines the requirement from physical orthodoxy that theories-of-consciousness must
65  obey *locality in physical causation*, and the consequent implications for brain-dynamical encoding of
66  consciousness. Although exceptions to locality within orthodox physical theory are sometimes
67  argued to occur in certain quantum-theoretic contexts, what is meant here by a "physically orthodox
68  theory-of-consciousness" (Baars, 1988; Edelman, 1989; Crick and Koch, 1990; Tononi and Edelman,
69  1998; Churchland, 2005; Baars and Edelman, 2012) is a member of the class of mainstream
70  explanations that are all based solely on classical physics. Thus, physically-orthodox theories-of-
71  consciousness must be local, in present nomenclature.

72  For an introduction to locality issues (**Figure 1**), consider as an illustrative example the encoding
73  of color at a particular location in the visual field, by the vector formed from the firing rates of three
74  different neurons (Solomon and Lennie, 2007). One typical mainstream proposal (Churchland, 1985)
75  is that the firing-rate vector is *identical to* or *synonymous with* the associated conscious experience of
76  color. Although a hypothesis of this kind can be advanced, it is *not* physically orthodox, because it
77  fails to be local in physical causation: the color-experience depends on information at three different
78  spatial locations. This simple observation perhaps evokes a jump directly to the central result, that
79  physically-orthodox (and therefore locally causal) theories-of-consciousness can only be valid if each
80  piece of fundamental information employed in the generation of conscious experience is encoded by
81  the brain in a spatiotemporally point-like manner. Clearly, however, such a leap would be premature,
82  because the example fails to exhaust all possibilities for physically orthodox theories-of-
83  consciousness. For example, consider the proposal that one of the emergent properties of the brain is
84  a spatial coordinate system in which the three vector-components *can* be considered co-located.
85  [Although this particular idea might appear esoteric, a complete and explicit account of space as it
86  relates to conscious experience (Dennett, 1991; Dennett and Kinsbourne, 1992) is required both

87    generally in the neuroscience of consciousness and specifically to treat locality properties rigorously.]
88    These considerations raise the primary question addressed in this paper: under what circumstances
89    can generation of consciousness from spatially-distributed codes satisfy the locality requirements of
90    physical orthodoxy?

91        A complete answer to this question must address the related contexts of *distributed* neural codes
92    and *parallel computation* (Rumelhart and McClelland, 1988; Thorpe and Imbert, 1989), because both
93    distributed codes and parallel computation employ physical dynamics at spatially-separated locations
94    to encode related information (**Figure 2**). In the present paper, a *distributed code* utilizes neural
95    dynamics in a number of spatially distinct neurons to encode a *single* piece of fundamental
96    information (for example, the specific orientation of an edge, or the specific hue of color, at a specific
97    external location). (Conventionally, the terminology "local code" is used to describe the
98    complementary case of single-neuron encoding. Because the term "local" will be used in this paper
99    solely as a descriptor of physical causality, instead the nomenclature "non-distributed" or "spatially-
100   point-like" will be used to describe single-neuron codes.) In contrast to the distributed/non-
101   distributed distinction that refers to basic encoding of a *single* fundamental piece of information, the
102   *parallel computation* nomenclature will be used here to refer to multiple distinct streams that encode
103   or process *multiple* pieces of related information (for example, edge-orientations at a number of
104   locations, or edge-orientation, color, and motion at a single location). Thus, in present terminology,
105   neural codes can in principle be non-distributed and non-parallel, distributed and parallel, non-
106   distributed and parallel, or distributed and parallel, although the human sensory encoding and
107   information processing that will form the focus of central developments is understood (Kaas, 1997;
108   Decharms and Zador, 2000) to be implemented, at least in early stages, via massively parallel
109   computations that typically employ distributed codes.

110       Data encoded at different spatial locations in the brain must be spatially integrated at some point
111   in processing, if physical action based on that information is to be locally causal. (Otherwise, brain
112   dynamics at multiple distinct locations would have a physical effect on some non-connected spatial
113   location, requiring definitively *non*-local causality.) Problems involved with the integration of
114   spatially-distributed encoding have been discussed to some extent under the heading of "the binding
115   problem" (von der Malsburg, 1981; Treisman, 1999; Engel and Singer, 2001; Di Lollo, 2012;
116   Feldman, 2013), although these literatures have failed to fully address significant aspects relevant to
117   the viability of physically-orthodox theories-of-*consciousness*, for four reasons. First, discussions of
118   the binding problem often fail to distinguish clearly between *the binding problem for behavioral*
119   *encoding* (in which spatially distributed information in parallel streams must converge in an orderly
120   yet flexible manner, to support neural computations governing behavior), and *the binding problem*
121   *for consciousness* (in which spatially-distributed information in parallel streams must be brought
122   together via some as-yet unknown mechanism, to generate the components of conscious experience
123   at each experiential location). Second, even those discussions that disambiguate and then focus on the
124   binding problem for consciousness fail to address *the relationship between spaces*, namely the
125   orthodox physical space (*i.e.* the unseen but putatively existent space in which elementary particles
126   exist) and conscious-experiential space (*e.g.* the visually experienced "emptiness" between distinct
127   objects). Discussion of these two spaces is essential to questions of local causation in the generation
128   of conscious experience, because the two spaces might themselves be topologically unconnected in
129   certain theories-of-consciousness (**Figure 3**), and thus, in a sense to be examined more closely in
130   Section 6, non-local to each other in a basic manner. Third, the question of *temporal locality* in
131   physical causation has been completely neglected, in the context of explaining the neural bases of
132   consciousness. Just as spatial locality becomes questionable when codes are extended across a

Nicholas M. Rosseinsky

133   number of spatial locations, temporal locality must be explicitly accounted for when codes are
134   extended across time in some manner. Fourth, examinations of the binding problem (and
135   consciousness more generally) are usually *couched in informal, verbal, terms*, and thus lack the
136   formal rigor necessary to properly address the behavioral/consciousness binding-problem distinction,
137   the physical-space/conscious-experiential-space relationship, or even a sufficiently complete
138   definition of spatiotemporal locality.

139   Discussions thus far indicate that a complete answer to questions of locality in the generation of
140   consciousness will need to address all possible physically-orthodox theories-of-consciousness, to
141   consider both distributed and parallel processing, and to remedy the noted deficits in prior
142   approaches. These demands are considerable, and evoke a formal symbolism (summarized in Table
143   1) for describing relevant aspects of environment, brain, and experience, that will bestow benefits of
144   both generality and precision. Foundations for this formalism are constructed in Section 2 via
145   orienting definitions of salient consciousness-related phenomena. Sections 3 and 4 develop formal
146   symbolism for describing, respectively, the external physical *environment* and consequent brain-
147   dynamical *encodings* of environmental features. Section 5 develops analogous formalism for
148   *conscious experience* of the environment, notably including a coordinate system for relative spatial
149   location of the components-of-consciousness. This coordinate system is then used in Section 6 to
150   describe relevant scenarios for the physical-space/conscious-experiential-space relationship
151   mentioned above, thus providing an informal introduction to locality issues. The relatively lengthy
152   investment in formalism made in Sections 2-5 yields benefits in Section 7, via succinct formal
153   demonstrations that spatially-distributed and temporally-extended codes cannot be the final
154   informational basis for consciousness in a physically orthodox setting. Section 8 describes how the
155   exclusion of these codes constitutes an empirical test for physically-orthodox theories-of-
156   consciousness, and discusses the theoretical implications of alternative empirical outcomes.

157

## 2.    Consciousness
### 2.1.    Definition

160   "Consciousness" can be defined *e.g.* as the collection of phenomena that are subjectively present in
161   the awake state of a typical human being, and that are in contrast absent during deep sleep. Although
162   an appeal to subjectivity appears problematic for scientific enquiry that purports objectivity, the
163   subjective/objective distinction must be scientifically acceptable: it is foundationally necessary for
164   the existence of science itself, precisely in order to delineate objective scientific constructions from
165   subjective quotidian experience (for example, to delineate the objective entities constituting the
166   physical environment from the subjective human experience of that environment). Serious problems
167   in the scientific consideration of consciousness *do* arise, however (Hawking, 2000), when theories
168   hypothesizing necessary and sufficient physical conditions for the generation of subjective
169   experience are stated, because there is currently no experimental method for establishing whether an
170   allegedly conscious artefact (Penrose, 1989; Tononi, 2004) does in fact possess subjective experience
171   of any kind. In the case of *human* consciousness, this problem can be partially ameliorated by the
172   presumption that each "typical" human brain has a (currently unknown) mechanism for generating
173   conscious experience, just as each typical human brain has mechanisms *e.g.* for processing afferent
174   sensory data in order to navigate the environment in biologically advantageous ways. In principle,
175   discovery of biophysical mechanisms critically involved in the generation of consciousness (but not
176   in other biological functions) might be experimentally verified in the future by techniques that

177 temporarily interfere with these mechanisms, which should result in disruption of associated
178 components of consciousness (but not of behavior).

179     The present paper avoids absence-of-experimental-technology problems associated with *general*
180 theories-of-consciousness, because it does not attempt to identify a unique, universally applicable,
181 statement of the physics of conscious experience. Instead, it accepts consciousness as a phenomenon
182 associated with the *human* brain (without any assessment at all of the conscious state of any other
183 biological or non-biological system), and then derives requirements for biophysical mechanisms that
184 must apply if theories-of-consciousness are to be of an entirely physically-orthodox kind. This
185 approach differs radically from most science-of-consciousness papers, whose goal is to establish
186 particular theory-of-consciousness: here, the goal is to establish empirical tests that delineate the
187 entire class of physically-orthodox, mainstream theories [5.2] from those theories requiring either
188 new physics or the unorthodox interpretation of orthodox physics. Thus, the present approach does
189 not depend on experimental tests for conscious *vs.* non-conscious systems, because theoretical
190 developments here do not make predictions concerning this distinction. Nevertheless, the approach
191 here is definitively scientific, rather than philosophical, because it leads to concrete empirical
192 determinations that can arbitrate between two significant theory-classes.

193

194 **2.2.    Limitation to exteroceptive consciousness**

195 Even with a common orienting definition [2.1], the terms "consciousness" and "conscious" are used
196 with a variety of different connotations, leading to the possibility of semantic misunderstanding.
197 Notably, "consciousness" is sometimes used to refer to the controversial idea of a causally-
198 efficacious subject (in statements such as "the 'conscious self' made this choice"), whereas present
199 considerations will explicitly *exclude* the conscious self, and instead focus on *exteroceptive*
200 consciousness, *i.e.* on experience of the external environment. That is, without any axiomatic
201 assumptions concerning the objective existence (Popper and Eccles, 1977; Hameroff, 2012) or
202 otherwise of the 'self', considerations here will be limited to the subset of consciousness-related
203 phenomena whose biophysical origination depends critically on sensory processing of objective
204 physical features in the external environment. Specifically, developments will be limited to visual
205 and auditory modalities (**Figure 4**), because this limitation of scope offers greatest expository clarity,
206 without loss of generality in results.

207

208 **2.3.    Axiomatic role for brain dynamics in providing information to consciousness**

209 "Psychophysical parallelism" proposes that correspondence between exteroceptive conscious
210 experience and the objective contents of the external environment need not depend on brain
211 dynamics. However, axiomatically in the present work, it will be assumed that brain dynamics are the
212 physical mechanism by which information concerning the environment is transferred to conscious
213 experience. Thus, brain dynamics will be variously described as "constituting the informational basis
214 for consciousness" or "generating a contribution to [or component of] consciousness", language
215 which explicitly acknowledges a causal role for brain dynamics without proposing any specific
216 theory-of-consciousness. Avoiding adherence to particular theories-of-consciousness is required by
217 the basic approach here of comparing and contrasting two large *classes* of theory [2.1].

218

## 3.    Formal symbolism for features of the exteroceptive environment

### 3.1.    $s_{a,b}$ notation for perceptual feature types and instances

The brain is thought to parse the external environment with respect to a small finite number of perceptual feature types for each modality (Kandel *et al*., 2012; Zaidi *et al*., 2013). For example, feature types associated with vision include color, motion, oriented edges, and so forth. For the auditory modality, encoding can be conceived in terms of a single feature type, *i.e.* amplitude, indexed secondarily by frequency. (That is, sound at a particular location consists of a set of amplitudes across a variety of frequencies). For each feature type, there are then a number of different specific instances that can be encoded by the brain. For example, for the color feature *type*, the brain is thought to be capable of distinguishing around $10^7$ different color *instances* (Judd and Wyszecki, 1975). For expository simplicity, it will be assumed that each feature type has only a finite number of discrete instances, although in principle a continuous code can generate an infinite number of instances. (Explicitly accommodating the continuous code possibility would considerably complicate symbolism, but not affect results.)

The notation $s_{a,b}$ will be adopted for the *b*-th instance of the *a*-th feature type (**Figure 5**). As a rudimentary illustrative example, if the first feature type $s_1$ is color in the visual field, then $s_{1,1}$ might denote red, $s_{1,2}$ orange, and so on. (More rigorously, $s_{1,1}$ and $s_{1,2}$ should be defined in terms of wavelengths of reflected light, in part because terms such as "red" refer ambiguously both to objective wavelengths and subjective experience: present notation aims to explicitly separate the objective and the subjective, so that the relationship between them can be treated clearly.) The adoption of the basic *s* notation is meant to reinforce that $s_{a,b}$ symbols label various *stimuli* in the exteroceptive environment.

Emphatically, $s_{a,b}$ symbols label *objective* stimulus features that can in principle be reduced to configurations of atoms and molecules, and thence to even more elementary physical constructs (quarks, electrons, and so on). For example, the perceptual feature type of visual color is a way of describing wavelengths of reflected photons, which could equally well be described in terms of the emission and absorption properties of atoms and molecules.

### 3.2.    Reserved usage of $\{\ldots\}$ notation

Parenthetic expressions of the form "{…}" (as opposed to "(…)" and "[…]") in the present paper explicitly connote multi-member sets or collections. Notably, later symbolism $A(\mathbf{r}_i)$ means the value of *A* at some *particular* coordinate location $\mathbf{r}_i$, whereas $\{A(\mathbf{r}_i)\}$ means the collection of *A*-values at *every* contextually-relevant $\mathbf{r}_i$. Similarly, in present notation, a function $f[A(\mathbf{r}_i)]$ depends only on the value of *A* at a single location, whereas a function $g[\{A(\mathbf{r}_i)\}]$ depends on *A* at multiple locations: this distinction between single- and multi-location dependence of physical action (or computation) will be pivotal in spatial-locality considerations.

### 3.3.    $_X\mathbf{r}_j$ notation for locations in the exteroceptive environment

257   Of course, navigation of the environment by the human organism requires not only identification of
258   perceptual features labelled by $\{s_{a,b}\}$, but also of the locations of these features relative to current
259   position. A complete formal description of the exteroceptive milieu for present purposes therefore
260   requires the assignment of coordinate indices to perceptual features. Because sensory organs (*e.g.* the
261   retina) consist of a finite number of cells with discrete receptive fields, it will be assumed that there
262   are a finite number of fixed environmental locations sampled in exteroception. (Various other
263   assumptions are possible, but do not affect results.) Adopting a conventional notation **r** to denote
264   coordinates with respect to some well-defined (*e.g.* brain-centred, Euclidean) spatial axes, the
265   coordinates of exteroceptively relevant locations will be denoted by $_X\mathbf{r}_j$ (**Figure 6**), where $j = 1, \ldots,$
266   $N_E$, and $N_E$ is the total number of locations. (The prefixed subscript "X" labels these coordinates as
267   "eXternal" locations at which various stimulus features $s_{a,b}$ exist, as opposed to later identifiers "A"
268   [4.4] and "B" [5.5.3] that denote two sets of internal brain locations at which dynamical activity has
269   encoding significance.) Thus, the set $\{_X\mathbf{r}_j\}$ comprises the coordinates of all locations in the
270   exteroceptive environment from which sensory information originates.

271

272   **3.4.    Completeness and generality of $\{s_{a,b}(_X\mathbf{r}_j)\}$ notation**

273   **3.4.1.  Completeness**

274   By construction, an explicit definition of the totality of $s_{a,b}$ at each and every $_X\mathbf{r}_j$ constitutes a
275   complete description of the exteroceptive environment, in the sense that neural responses in sensory
276   organs are fully entailed by this information (together with complete understanding of the human
277   organism's neural systems, and full knowledge of the non-sensory states of the organism, such as
278   wakefulness). Note that complete description of the visual environment typically requires the
279   specification of many $s_{a,b}$ at each non-empty $_X\mathbf{r}_j$, because, minimally, geometry (*e.g.* orientation of
280   edges), color, and motion must all be specified, and these correspond to three different feature or
281   stimulus types.

282

283   **3.4.2.  Generality**

284   The collection $\{s_{a,b}(_X\mathbf{r}_j)\}$ is simply a way of describing the external environment with respect to
285   particular set of categories that map in a very direct way to neural encoding (because symbols were
286   defined with this goal in mind). Note that symbols are theory-of-encoding independent, in the
287   following sense. Despite the fact that symbols were developed above with reference to well-
288   established roles for specific features in neural encoding (such as edges, color, and motion, for
289   vision), symbolically-stated theories in the present paper do not rely on the ultimate validity of these
290   roles. In the limit, $\{s_{a,b}(_X\mathbf{r}_j)\}$ symbols can be taken to stand for the totality of elementary physical
291   constituents of the exteroceptive environment, thus liberating analyses from dependence on any
292   particular theory of encoding, although this approach considerably complicates accompanying verbal
293   exposition (and will not be adopted here for that reason).

294

295   **4.      Formal symbolism for neural encoding of the exteroceptive environment**

Nicholas M. Rosseinsky

### 4.1. Theories of neural encoding

For the purposes of explaining behavior, neural dynamics in sensory organs can be conceived of as encoding the environment (Perkel and Bullock, 1968; Field, 1994; Rieke *et al.*, 1997), in the sense that perfect knowledge of sensory physiology can be used to make certain inferences concerning environmental contents, based on observations of neural dynamics. The relationship between sensory encoding and later motor activity can further be conceived of as a series of computational transformations of the initial encoding (Phillips *et al.*, 1984; Churchland and Sejnowski, 1992), although rigid analogies with human-fabricated electronic computers must be moderated *e.g.* by the role in the brain of contextual dynamical oscillations (Buzsáki, 2006) that have no direct computer analogue. Moreover, despite tremendous advances in the understanding of various types of encoding and computation in sensory and later neural processes [*e.g.* (Olshausen, 1996; Bell and Sejnowski, 1996; Rodriguez *et al.*, 1999; Abbott and Sejnowski, 1999; Mehta *et al.*, 2002; Haynes and Rees, 2006; Quiroga and Panzeri, 2009; Panzeri *et al.*, 2010; Horikawa *et al.*, 2013)], there are many outstanding areas of significant uncertainty (Eggermont, 1998; Van Vreeswijk, 2004) concerning, for example, fidelity of codes in the presence of "noisy" stochastic behavior (Mainen and Sejnowski, 1995; Faisal *et al.*, 2008), the basic organizational unit of encoding (Abeles, 2011), the binding of information encoded in multiple cortical areas (Feldman, 2013), and the encoding role of spatially-distributed electromagnetic-field oscillations (Buzsáki, 2006; Sejnowski and Paulsen, 2006).

In summary, although there is little doubt that the neural coding description is a helpful picture of brain function, currently there is no complete understanding of how information is represented and transformed in neural processing. Thus, any description of how brain encoding contributes information to conscious experience [2.3] must at present be developed in a relatively generic way, specifically via statements that avoid stating a *particular* theory of encoding.

### 4.2. Notation *A* for brain-dynamical encoding activity

Neural encoding of the exteroceptive environment means that certain brain-dynamical states, but not others, contain the information that a particular instance of a particular perceptual feature type is present at a particular environmental location. In order to formally state neural coding theories, it is necessary to specify a physical *measure* that can distinguish encoding and non-encoding states. The symbol *A* will be used to denote this measure.

### 4.3. What does *A* stand for?

Three empirical facts support the frequent adoption of an electromagnetic field ("e.m.-field") measure as means of defining and delineating dynamical states relevant to behavioral encoding. First, initial sensory encoding takes place via disturbance of membrane potentials. Second, propagation of information is primarily via electrochemical (axon-to-dendrite-to-axon) and electrical (gap junction) mechanisms. Third, behavioral computations are ultimately converted to motor signals of an electrochemical form. Naturally, it is possible to restate the electrical state of the brain in chemical or biochemical terms, for example, characterizing a "firing" neuron by ion flows or by channel states, rather than by membrane potential. However, for present purposes, these different descriptions are ultimately equivalent, so it will be assumed for discursive simplicity that the behavioral encoding

337  state of the brain is defined by e.m.-field dynamics. Thus, $A$ can be thought of as standing for some
338  measure of e.m.-field state, such as membrane potential.

339

340  **4.4.    $_A\mathbf{r}_i$ notation for encoding-relevant brain locations**

341  To understand what a brain is encoding, one must know the value of $A$ at a set of locations whose
342  coordinates will be denoted by the generic symbol $_A\mathbf{r}_i$ (**Figure 7**). (That is, the set $\{_A\mathbf{r}_i\}$ comprises the
343  coordinates of all brain locations at which $A$ must be known in order to deduce encoded features of
344  the exteroceptive environment. As mentioned in [3.3], the prefixed "A" subscript distinguishes *e.g.*
345  $_A\mathbf{r}_1$ from $_X\mathbf{r}_1$: the former indexes a brain-encoding site, the latter an external environmental location.)
346  The precise loci labeled by the $\{_A\mathbf{r}_i\}$ depend on how the human brain actually encodes exteroception
347  [4.1], but this affects neither their generic definition, nor present developments. For example, the
348  collection $\{_A\mathbf{r}_i\}$ can index *every* brain location.

349

350  **4.5.    $C_{abj}$ notation for activity engendered by $s_{a,b}(_X\mathbf{r}_j)$**

351  By construction, the collection of $A$ measurements at the locations $\{_A\mathbf{r}_i\}$ are sufficient to define the
352  encoding state of the brain. Given a specific stimulus $s_{1,2}$ and a specific environmental location $_X\mathbf{r}_3$,
353  say, it is therefore possible to parse the information $\{A(_A\mathbf{r}_i)\}$ to determine whether $s_{1,2}(_X\mathbf{r}_3)$ is
354  encoded. The precise details of parsing depend again on how the brain actually encodes the
355  environment [4.1], but any reliable encoding method must create a partition between those $A$-
356  measured brain states that encode a particular $s_{a,b}(_X\mathbf{r}_j)$ (in the general case) and those that do not. This
357  fact can be formally encapsulated by introducing a classifier function $C_{abj}$ on the collection $\{A(_A\mathbf{r}_i)\}$,
358  defined via

359  $$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1 \iff \{A(_A\mathbf{r}_i)\} \text{ encodes } s_{a,b}(_X\mathbf{r}_j) \tag{1}$$

360  and

361  $$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 0 \iff \{A(_A\mathbf{r}_i)\} \text{ does not encode } s_{a,b}(_X\mathbf{r}_j) \tag{2}$$

362  (where $\iff$ means "if and only if"). The function $C_{abj}$ is then a generic formal placeholder for any
363  specific theory of encoding. For example, if $s_{a,b}(_X\mathbf{r}_j)$ is encoded by the membrane potential of a single
364  neuron, then $C_{abj}$ simply corresponds to the assessment of potential relative to threshold [4.6]. On the
365  other hand, in a complex spatiotemporal encoding, $C_{abj}$ corresponds to determining the presence or
366  otherwise of the particular distributed pattern (or patterns) that encode $s_{a,b}(_X\mathbf{r}_j)$. (Note that, as a
367  default, the argument of $C_{abj}[\ldots]$ is defined as the *multi*-location collection $\{A(_A\mathbf{r}_i)\}$ of $A$-values
368  [3.2].)

369      Let the symbols $\exists_P$ and $\nexists_P$ respectively denote the physical (subscript "P") existence and non-
370  existence of an object in a particular setting. Specifically, $\exists_P s_{a,b}(_X\mathbf{r}_j)$ will mean that the stimulus $s_{a,b}$
371  exists in the physical environment at $_X\mathbf{r}_j$, whereas $\nexists_P s_{a,b}(_X\mathbf{r}_j)$ indicates that $s_{a,b}$ does not exist at that
372  location. Then, if encoding is a faithful representation of the environment, it follows that

373  $$\exists_P s_{a,b}(_X\mathbf{r}_j) \implies C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1 \tag{3}$$

374    and

375         $\nexists_P s_{a,b}(_X\mathbf{r}_j) \Rightarrow C_{abj}[\{A(_A\mathbf{r}_i)\}] = 0$                              (4)

376    (where $\Rightarrow$ denotes "… implies that …").

377

378    **4.6.      Example: $C_{abj}$ notation for a single neuron, membrane potential, encoding scheme**

379    A concrete example may help to clarify notation introduced in [4.5]. Consider the stimulus $s_{1,2}$ at the
380    location $_X\mathbf{r}_3$, and a hypothetical neural encoding in which the presence or absence of $s_{1,2}(_X\mathbf{r}_3)$ leads
381    respectively to firing or non-firing of a neuron whose axon hillock is labelled by coordinates $_A\mathbf{r}_4$. If
382    $A_{\text{thresh}}$ is the minimal membrane potential that is always and only exceeded during firing, then the
383    stimulus-firing relationship can be formally stated as

384         $\exists_P s_{1,2}(_X\mathbf{r}_3) \Rightarrow A(_A\mathbf{r}_4) > A_{\text{thresh}}$                              (5)

385    and

386         $\nexists_P s_{1,2}(_X\mathbf{r}_3) \Rightarrow A(_A\mathbf{r}_4) \leq A_{\text{thresh}}$                             (6).

387    Eqs. 5-6 simply state the physical and biophysical logic that presence or absence of stimulus leads to
388    firing or non-firing. Recalling that $C_{123}$ is defined (Eqs. 1-2) with respect to whether $\{A(_A\mathbf{r}_i)\}$ encodes
389    $s_{1,2}(_X\mathbf{r}_3)$, it follows that

390         $C_{123}[\{A(_A\mathbf{r}_i)\}] = 1 \iff A(_A\mathbf{r}_4) > A_{\text{thresh}}$                         (7)

391    and

392         $C_{123}[\{A(_A\mathbf{r}_i)\}] = 0 \iff A(_A\mathbf{r}_4) \leq A_{\text{thresh}}$                        (8).

393    Note that, although $C_{123}$ on the left-hand side ("LHS") of Eqs. 7-8 takes the complete set of values
394    $\{A(_A\mathbf{r}_i)\}$ as its argument, its one-or-zero value is set only with reference to $A(_A\mathbf{r}_4)$. Thus, the large set
395    of all possible combinatorial instances of $\{A(_A\mathbf{r}_i)\}$ is partitioned into the subset in which $A(_A\mathbf{r}_4) >$
396    $A_{\text{thresh}}$ (independent of $A$-values at other $_A\mathbf{r}_i$) and the disjoint subset in which $A(_A\mathbf{r}_4) \leq A_{\text{thresh}}$ (again
397    independently of other $A$-values).

398        $C_{123}$ can then be used to replace the $A(_A\mathbf{r}_4)$ firing/non-firing inequalities on the right-hand side
399    ("RHS") of Eqs. 5-6, to give

400         $\exists_P s_{1,2}(_X\mathbf{r}_3) \Rightarrow C_{123}[\{A(_A\mathbf{r}_i)\}] = 1$                          (9)

401    and

402         $\nexists_P s_{1,2}(_X\mathbf{r}_3) \Rightarrow C_{123}[\{A(_A\mathbf{r}_i)\}] = 0$                         (10).

403    For this particular encoding, Eqs. 9 and 10 are then directly the $a = 1$, $b = 2$, $j = 3$ instances of the $C_{abj}$
404    relationships in Eqs. 3 and 4.

405

### 4.7. Notation capable of describing any encoding scheme

Eqs. 3-4 are, by construction, simply very general formal statements of the basic premise of any theory of neural encoding [4.1], namely that the presence or absence of a particular stimulus in the environment leads to distinguishable patterns of brain activity. Thus, notation is clearly capable of describing any neural encoding scheme.

411

## 5. Formal symbolism for conscious experience of the exteroceptive environment

### 5.1. Approach to theories-of-consciousness

As noted in [2.1], the present paper does not aim to state or prove a particular theory-of-consciousness. Rather, it discusses the consequences of locality in physical causation for *all* theories of consciousness (after the exclusion of psychophysical parallelism [2.3]). In particular, results here will distinguish between biophysical encoding requirements for "physically orthodox" theories of consciousness [5.2], and for the complementary collection of theories that do not meet all the defining requirements of physical orthodoxy.

420

### 5.2. "Physically orthodox" theories-of-consciousness

As noted in the Introduction, the current mainstream assumption in physics and neuroscience is that consciousness can be explained within a physically-orthodox setting, either as synonymous with, or as an emergent property of, brain dynamics. Notably, the mainstream view denies the need for "new physics" of any kind in the explanation of consciousness. A "physically-orthodox setting" is taken here to *minimally* possess three specific properties, namely that physical causation is local, that there are only three spatial dimensions, and that physical structures are composed only of Standard Model elementary particles, such as quarks, electrons, and photons. Although the second and third assumptions are undoubtedly implicit in every mainstream theory-of-consciousness, their explicit emphasis might appear somewhat arbitrary, or even peculiar. However, as described in [7.2.6], these two uncontroversial assumptions (in mainstream settings) create significant limits on *how* the contents of conscious experience can be physically constituted relative to brain dynamics, by ruling out intermediary connective structure that might in principle assimilate distributed codes, and thus overcome locality problems.

435

### 5.3. $\langle s_{a,b} \rangle$ notation for conscious experience of perceptual feature types and instances

Formal symbolism thus far can describe the external environment [via symbols $s_{a,b}(_X\mathbf{r}_j)$], brain activity [via symbols $A(_A\mathbf{r}_i)$], and the relationship between environment and brain activity [via symbols $C_{abj}$]. However, no symbolism has yet been declared for *conscious experience* of the environment. The definition of $s_{a,b}$ symbols allows a very direct construction of required symbolism for consciousness. Introducing notation $\langle ... \rangle$ to denote "the conscious experience that results from the existence of ... in the exteroceptive environment", the total conscious exteroceptive experience

443  can be written as $<\{s_{a,b}(\mathrm{x}\mathbf{r}_j)\}>$. Similarly, the contribution to conscious experience from a single
444  stimulus $s_{a,b}(\mathrm{x}\mathbf{r}_j)$ is formally labelled as $< s_{a,b}(\mathrm{x}\mathbf{r}_j) >$ (**Figure 8**).

445

## 5.4.    $\rho_k$ notation for locations within exteroceptive conscious experience

447  Self-evidently [2.1], exteroceptive experience consists of a variety of stimulus-induced components
448  arranged in an orderly manner in a spatially-extended setting. The generic symbol $\rho$ will be used to
449  denote a location within the space of conscious experience (**Figure 9**), just as the symbol $\mathbf{r}$ [3.3] is
450  used to index objective space. As a simplification (that does not affect results), it will be assumed
451  that exteroceptive conscious experience can be fully described by specifying the contents of
452  consciousness at a finite number of locations denoted by $\{\rho_k\}$. Specifically, if sensory information is
453  completely defined by specifying the set $\{s_{a,b}\}$ at $\mathrm{x}\mathbf{r}_i$, where $i = 1, \ldots, N_E$ [3.3], then exteroceptive
454  conscious experience is completely defined by specifying the $\{<s_{a,b}>\}$ at each $\rho_k$, where $k = 1, \ldots, N_E$
455  and $\rho_k$ labels the location in conscious-experiential space corresponding to $\mathrm{x}\mathbf{r}_k$ in objective space.

456      The introduction of spatial coordinates for conscious experience raises the typically neglected
457  issue of the relationship between objective and subjective spaces (Dennett and Kinsbourne, 1992).
458  For example, consider the "space" comprised of the emptiness between and surrounding various
459  well-defined objects in visual conscious experience. What, if any, is the correlate in objective reality
460  of experiential space? This complex question cannot be answered in any scientifically definitive way
461  at present. (For example, unobservable physically-orthodox "space" only has a concrete existence in
462  a conventionally realist ontology, whose unique validity cannot currently be established.) For present
463  purposes it is only necessary to establish two points. First, the attribution of coordinates to locations
464  in conscious experiential space is scientifically valid [5.4.1]. Second, the introduction of the generic
465  and specific symbols $\rho$ and $\rho_k$ is entirely general (*i.e.* does not introduce any particular theory-of-
466  reality or theory-of-consciousness) [5.4.2].

467

### 5.4.1.   Scientific validity of coordinates for locations in conscious experiential space

469  Although not typically considered explicitly by computational neuroscience, $\rho$ coordinates (or their
470  informal proxies) have been long studied within visual psychophysics (Luneberg, 1944; Foley, 1978;
471  Heller, 1997), where one significant focus has been the characterization of visual as opposed to
472  objective geometry (Wagner, 2006). Consider, for example, the subjective visual experience of the
473  sky on a clear night, in which far distant stars are apparently embedded in a spherical dome,
474  contrasting with the approximately flat spatial geometry understood (*e.g.* by relativistic cosmology)
475  to pertain objectively. This simple comparison shows that subjective and objective geometries need
476  not coincide, and has led to extensive and sophisticated investigations of the general relationship
477  between the two [*e.g.* (Foley *et al.*, 2004)]. In this context, the $\rho$ coordinate system is the basic
478  labelling nomenclature for relative locations within subjective geometry, whereas $\mathbf{r}$ coordinates are
479  the basis of objective geometry.

480

### 5.4.2.   Generality of $\{\mathbf{r}, \rho\}$ formal symbolism

482  The fact of two geometries [5.4.1] does not at all *imply* physical duality of spaces, *i.e.* that locations
483  labeled by $\rho$ coordinates and those labeled by **r** coordinates constitute two ontologically equivalent
484  but physically unconnected domains (although neither is this possibility *excluded*, in the general
485  case). For example, the space comprised of locations labeled by $\rho$ coordinates might, as a whole, be a
486  *property of* objective dynamics occurring in the space of locations labeled by **r** coordinates, so that
487  the two spaces are physically and ontologically incommensurable. Another possibility is that objects
488  with  $\rho$ coordinates and those with **r** coordinates might exist in the *same* physical space, so that $\rho$
489  labels are just different names for points already conventionally labeled by **r** coordinates, thus
490  completely eliminating two-ness of constructs. The point for present purposes is that these
491  possibilities (and others) are simultaneously accommodated by the scientifically valid [5.4.1]
492  introduction of $\rho$ coordinates *without* accompanying reduction in current uncertainties concerning the
493  physical, ontological and topological aspects of the $\rho$-**r** relationship (other than the required
494  establishment of a *representational* relationship between $\rho_k$ and $\mathbf{r}_k$, namely that $\{s_{a,b}\}$ stimulus
495  features at $\mathbf{r}_k$ are represented by $\{<s_{a,b}>\}$ experiences at $\rho_k$).

496

497  ### 5.5.    *B* notation for final brain-dynamical cause of contribution to conscious experience

498  ### 5.5.1.  Naïve theory-of-consciousness employing $C_{abj}[\{A(_A\mathbf{r}_i)\}]$

499  It is perhaps tempting to move directly from notation developed so far to a formal statement of a
500  general theory-of-consciousness, as follows. Because the presence of $s_{a,b}(_X\mathbf{r}_j)$  in the exteroceptive
501  environment causes brain dynamics satisfying $C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$, it seems reasonable to propose that
502  $C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$ is a necessary and sufficient condition for the existence of $<s_{a,b}>(\rho_j)$  (where
503  $<s_{a,b}>(\rho_j)$ denotes the existence of $<s_{a,b}>$ at $\rho_j$ [5.4]), *i.e.*

504  $$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1 \Leftrightarrow \exists_P <s_{a,b}>(\rho_j) \tag{11}$$

505  and

506  $$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 0 \Leftrightarrow \nexists_P <s_{a,b}>(\rho_j) \tag{12}.$$

507  Notably, because existence of $s_{a,b}$ at $_X\mathbf{r}_j$ implies that $C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$ (Eq. 3), Eq. 11 implies, as
508  expected, that

509  $$\exists_P s_{a,b}(_X\mathbf{r}_j) \Rightarrow \exists_P <s_{a,b}>(\rho_j) \tag{13}$$

510  *i.e.* the existence of a stimulus in the exteroceptive environment leads to the corresponding conscious
511  experience at the appropriate location in subjective space.

512

513  ### 5.5.2.  Biomolecule basis-of-consciousness example: failure of $C_{abj}[\{A(_A\mathbf{r}_i)\}]$ theories

514  Although Eqs. 11-12 are a significant advance (offering the first formal expression of a relationship
515  between brain dynamics and components of conscious experience), they overstate causality and fail
516  to provide an entirely general theory-of-consciousness, as demonstrated by the following example.

Nicholas M. Rosseinsky

517       Consider a hypothetical situation in which the contribution-to-consciousness from the dynamical
518   activity of a particular neuron is governed, as a fact of nature, by the state of a specific biomolecule
519   $M$, rather than by the values of the e.m.-field measure $A$. Because $A$ is presumed to be the general
520   carrier of information [4.3], the consciousness-relevant aspect of $M$ must couple physically with $A$ in
521   a reasonably close way, so that *e.g.* neuronal firing always puts $M$ into the state (or states) that
522   generate the appropriate component of conscious experience. For example, if $M$ has a non-uniform
523   charge distribution (so that e.m-field disturbances from firing can affect $M$-state), firing can cause a
524   particular, hypothetically consciousness-relevant, energy state of $M$. The point is that, in this
525   example, it is some aspect of $M$ that is the final physical cause of contribution-to-consciousness, not
526   simply the $A$-state of the neuron.

527       When applied to this example, Eqs. 11-12 *overstate causality*, because $C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$ is not
528   sufficient to ensure the generation of the conscious experience $<s_{a,b}>(\rho_j)$: if the $M$ molecule is
529   missing or defective, the neuron can be in an encoding $A$-state but no contribution to consciousness
530   will occur. Formalism thus far is also *an incomplete account of causality* for this example, because
531   there are no formal symbols yet for the causal role of $M$.

532

### 5.5.3. Definition of $B$ and $_B\mathbf{r}_i$

534   The precise final physical cause of brain-dynamical contribution-to-consciousness is not presently
535   known: it could be the $A$-state of a neuron, the state of some dedicated $A$-coupled biomolecule $M$
536   [5.5.2], or any of a number of other physical features that are coupled to or accompany $A$-dynamics
537   (*e.g.* ion flows or concentrations, vesicle exocytosis, channel states, and so on.) Thus, any general
538   formalism for theories-of-consciousness must accommodate the possibility of a final physical cause
539   for contribution-to-consciousness other than $A$. The symbol $B$ will be used denote a measure of this
540   *final physical cause* [with the understanding that, in e.m.-field-based theories-of-consciousness
541   (McFadden, 2002; Pockett, 2002), $B$ and $A$ are two symbols for one physical measure]. For example,
542   if energy of a biomolecule is the final physical cause (as imagined in [5.5.2]), $B$ might be defined in
543   terms of a set of distances relative to fixed, local, cellular structure.

544       Because the location at which $B$ is measured need not be the same as that at which $A$ is measured,
545   the symbols $\{_B\mathbf{r}_i\}$ will be introduced to denote the $B$-measurement locations. For example, in the
546   biomolecule case [5.5.2], $_A\mathbf{r}_i$ might denote the location of the axon hillock, and $_B\mathbf{r}_i$ the (nearby)
547   location of the biomolecule $M$.

548

### 5.6. $D_{abj}$ notation for contribution to conscious experience engendered by $s_{a,b}(_X\mathbf{r}_j)$

550   Just as $C_{abj}$ classifies those states of $A$ that encode $s_{a,b}(_X\mathbf{r}_j)$, so the function $D_{abj}$ is defined to classify
551   those $B$-states that generate $<s_{a,b}>(\rho_j)$. Specifically, as direct analogies to Eqs. 1-2,

552 $$D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1 \iff <s_{a,b}>(\rho_j) \text{ exists in conscious experience} \qquad (14)$$

553   and

554 $\qquad D_{abj}[\{B(_B\mathbf{r}_i)\}] = 0 \Leftrightarrow <s_{a,b}>(\rho_j)$ does not exist in conscious experience $\qquad$ (15).

555 To complete a theory-of-consciousness (**Figure 10**), a physical coupling [5.5.2] between *A*-states and
556 *B*-states must be declared, *e.g.* via

557 $\qquad C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1 \Leftrightarrow D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1 \qquad$ (16).

558 After declaration of Eqs. 14-16, the causal chain from stimulus to conscious experience (Eq. 13) still
559 holds, but now depends on a $C_{abj}$-to-$D_{abj}$-to-$<s_{a,b}>(\rho_j)$ path, thus remedying the two deficits noted in
560 [5.5.2] for a purely $C_{abj}$ theory.

561

## 562 5.7.  Notation capable of describing any theory-of-consciousness

563 Notation so far does not limit subsequent developments to any particular theory of consciousness
564 (other than requiring that some *physical B*-measure participate in specification of the contents of
565 consciousness [2.3]). That is, $D_{abj}[\{B(_B\mathbf{r}_i)\}]$ formalism (Eqs. 14-16) is consistent with a wide range of
566 proposals, notably including all physically-orthodox theories. Further development of notation is
567 certainly required to delineate *formally* between various theoretical sub-classes, for example, to
568 distinguish between identity theories (Place, 1956; Smart, 1959) and emergent-property theories
569 (Van Gulick, 2001). These sub-classes make distinct proposals concerning the precise physical
570 properties that $D_{abj}[\{B(_B\mathbf{r}_i)\}]$ conditions modulate in order to generate conscious experience.
571 However, formalism for these distinctions will not be developed here, because symbolism for
572 coordinate locations and for causal effects already explicitly present in Eqs. 14-16 fully suffices for
573 locality considerations.

574

## 575 6.  The relationship between r-space and ρ-space

576 Notational developments thus far (summarized in Table 1) have been entirely general. (Although
577 certain defining aspects of physically-orthodox theories-of-consciousness were noted in [5.2], these
578 limitations were not used to constrain subsequent formalism.) In this Section and the next, the
579 approach shifts to consider only physically-orthodox theories, thus dividing theories into two classes.
580 The present Section further divides physically-orthodox theories into two sub-classes, according to
581 their treatment of the relationship between **r**-space and **ρ**-space [5.4]. This sub-division of orthodox
582 theories clarifies the basis for assessing locality.

583

## 584 6.1.  Physical orthodoxy admits only two basic conceptions of ρ-space

585 As noted in [5.4.2], **ρ**-**r** symbolism is a general notation, and can describe both dualist and non-
586 dualist theories-of-consciousness. However, under the limitation to physical orthodoxy [5.2], there
587 can be no physically-dual space. Moreover, a physically orthodox approach rules out higher-
588 dimensional constructions that embed orthodox three-space and conscious-space as non-overlapping
589 (disjoint, orthogonal) subspaces of a larger, topologically-connected space (Rosseinsky, 2014b). This
590 means that only two logical possibilities exist for the relationship between **r**-space and **ρ**-space. In

Nicholas M. Rosseinsky

591  the first possibility (here termed the "sub-domain" approach), ρ-coordinates label a sub-domain of
592  orthodox three-dimensional space that is employed for the generation of consciousness (**Figure
593  11A,B**). In the second (the "property" approach), the conscious-experiential space itself is a property
594  of brain dynamics (**Figure 11C**), so that ρ-coordinates label a kind of virtual or contingent "mental
595  space". This second case contains two distinct spatial constructs, namely the orthodox three-
596  dimensional space and the space-as-property, and may therefore appear to be a dualist conception
597  antithetical to physical orthodoxy. However, it is possible to construe space-as-property as having an
598  ontologically junior status, thus in a sense preserving the "only three [ontologically fundamental]
599  spatial dimensions" requirement of physical orthodoxy [5.2].

600

## 6.2. Locality in sub-domain and property conceptions of ρ-space

602  Physically-orthodox theories must employ either a sub-domain or a property approach to construct a
603  space *in which* $<s_{ab}>(\rho)$ experiences occur. This observation is useful to locality considerations,
604  because the two approaches raise different issues.

605    In the *sub-domain* approach [6.2.1], $<s_{ab}>(\rho)$ experiences and *B*-dynamics that encode them occur
606  in the same **r**-space, so that assessing locality is the rather straightforward matter of whether $<s_{ab}>(\rho)$
607  and *B*-dynamics are co-located. (For completeness, the discussion of this point in [6.2.1] will also
608  give an introductory treatment of the slightly less straightforward matter of how $<s_{ab}>(\rho)$ components
609  come to be arranged appropriately within conscious experience. This question arises because locality
610  constrains the **r**-location of $<s_{ab}>(\rho)$ to be the same as its encoding dynamic, thus arranging the $<s_{ab}>$
611  in **r**-space according to the relative locations of neurons.)

612    In the *property* approach [6.2.2], arrangement of $<s_{ab}>(\rho)$ components within conscious
613  experience is not problematic (because the existence of $<s_{ab}>$ at $\rho_j$ is taken to be a property of certain
614  *B*-states, and "at $\rho_j$" in a property approach means "at the location indexed by $\rho_j$ within conscious-
615  experiential space *that is independent of orthodox space*"). In property-based theories the key issue is
616  instead the conditions under which an **r**-location and a ρ-location can *ever* be "local" to each other,
617  because the property conception inherently proposes that there is no meaningful distance metric
618  between **r**- and ρ-locations. (If such a metric existed, **r**- and ρ-spaces are embedded in a common
619  metric space. Although this is a possible theoretical proposal, it does not fall into the class of
620  physically orthodox theories because common embedding implies a definite ontological equivalence
621  between **r**- and ρ-spaces [6.1].)

622

### 6.2.1. Locality in the sub-domain approach

### 6.2.1.1. Non-distributed and distributed codes in the sub-domain approach

625  Let $\mathbf{r}(\rho_j)$ denote the **r** coordinate that can be assigned to any $\rho_j$ in the sub-domain approach (but not
626  the property approach). Then the construction of conscious experience is only physically local if the
627  **r**-coordinate of the *B*-dynamics generating $<s_{ab}>(\rho_j)$ is identical to $\mathbf{r}(\rho_j)$. Clearly, this means that all
628  *B*-dynamics generating $<s_{ab}>(\rho_j)$ must be at the *same* location [namely, $\mathbf{r}(\rho_j)$]. Because distributed
629  codes employ *B*-dynamics at multiple distinct locations, they cannot be local. Additionally, a non-

630    distributed code must create $<s_{ab}>(\rho_j)$ at the same **r**-location as its encoding *B*-dynamic, if the
631    construction of experience is to be local.

632

### 6.2.1.2.           The construction of subjective geometry in the sub-domain approach

634    A local sub-domain approach must employ a non-distributed code in which $<s_{ab}>(\rho_j)$ is generated at
635    the same **r**-location as its encoding *B*-dynamic [6.2.1.1]. Because the space of conscious experience
636    is orthodox physical space in a sub-domain approach, this means that the collection $\{<s_{ab}>(\rho_j)\}$ will
637    be arranged within conscious experience *just as* corresponding *B*-encodings are arranged within
638    orthodox physical space. There are then two possibilities. First, if *B*-encodings are arranged within
639    orthodox physical space in a manner that recapitulates the geometry of the $\{<s_{ab}>(\rho_j)\}$ within
640    subjective experience [5.4.1], then (at least for a single feature type) brain dynamics generate
641    conscious experience in a direct and consistent manner. [Additional complications for multiple
642    feature types are discussed in (Rosseinsky, 2014b).] Second, if the **r**-configuration of *B*-encodings
643    does *not* recapitulate the subjective geometry of the $\{<s_{ab}>(\rho_j)\}$, some further mechanism must exist
644    to rearrange the $\{<s_{ab}>(\rho_j)$. To facilitate this rearrangement, it is possible to propose that space exists
645    under two metrics, the conventional **r**-metric and a **ρ**-metric that puts points close together in
646    conscious experience that are far apart under the conventional metric. [Note that elaboration of this
647    "dual metric" approach is necessary for a rigorously comprehensive treatment of physically orthodox
648    theories-of-consciousness, in the following sense. In the absence of a method for generating
649    subjective geometry from arbitrarily arranged *B*-dynamics, local sub-domain theories can only be
650    given if the **r**-configuration of *B*-dynamics already recapitulates subjective geometry, thus severely
651    limiting the physically-orthodox class. The question of whether dual-metric theories themselves
652    qualify as "physically orthodox" is discussed elsewhere (Rosseinsky, 2014b). For present purposes,
653    the approach will be to establish the *broadest possible* set of physically orthodox theories, and then to
654    show that, within this set, only non-distributed codes are local.]

655

### 6.2.2.  Locality in the property approach

### 6.2.2.1.          r/ρ-locality

658    In one view, every property-based theory is *non-local*, because points indexed by **ρ**-coordinates have
659    no meaningful distance from any **r**-location. In this view, the inability to establish locality (*i.e.* a *zero*
660    distance) is non-locality, leading to simplification of the locality-of-physically-orthodox-theories
661    question: only sub-domain theories can be local (and therefore, after [6.2.1.1], the only local
662    physically-orthodox theories are non-distributed codes in which $<s_{ab}>(\rho_j)$ is generated at the same **r**-
663    location as its encoding *B*-dynamic).

664        In an alternate view, every theory employing the property approach is *local*, again because **ρ**-
665    coordinates have no meaningful distance from any **r**-location! In this view, the inability to establish
666    non-locality (*i.e.* a *non*-zero distance) is locality. One notable extension of this view is that a **ρ**-
667    location can then be local (or more precisely "not **r**-non-local") to a multiplicity of **r**-locations. Note
668    that *B*-properties are defined at **r**-locations rather **ρ**-locations. But *if B*-properties could exist "in" **ρ**-
669    space, it would then be possible for a single **ρ**-location to "contain" *B*-values from a number of **r**-

670   locations, so that a multiple-**r** (*i.e.* distributed) code could be local with respect to ρ. Allowing this
671   chain of speculative possibilities would then lead to the result that distributed codes *can be* local in
672   certain property approaches, an implication of fundamental significance to central considerations

673       These discussions establish that the way in which locality is attributed to property approaches
674   (**Figure 3D,E**) will profoundly affect the final view of locality in physically-orthodox theories.
675   Therefore, if final results are to be reliable, an objective way of making this attribution must be
676   established. This can be done in two steps, as follows. First, because the central result will *exclude*
677   distributed codes from local, physically-orthodox, theories-of-consciousness, the attribution should
678   be made in the first instance in a manner that *includes* as many distributed coding possibilities as
679   possible. This means adopting the view that property approaches are naturally local, or more
680   precisely "not **r**-non-local", by virtue of the metric absence. Second, the locality (or not-non-locality)
681   thus made available to property approaches must be limited by the overall conception of property
682   approaches as *physically orthodox*. Recalling that concerns of a non-orthodox duality in property
683   approaches are to be answered by an ontologically junior status for conscious-experiential space
684   [6.1], this means that *B*-dynamics cannot exist in ρ-space, resulting in a critical constraint on the
685   locality of property approaches: **r**/ρ-locality in a property approach cannot be used to facilitate or
686   perform joint computations across *B*-values at multiple **r**-locations.

687

### 6.2.2.2.       Non-distributed and distributed codes in the property approach

689   The definition of **r**/ρ-locality established in [6.2.2.1] has the following implications. First, *non*-
690   distributed codes in the property approach are all local. (Note that in the sub-domain approach, it is
691   possible for a *non*-distributed code to be *non*-local, if $\mathbf{r}(\rho_j)$ is not identical to the $_B\mathbf{r}$ coordinate of the
692   corresponding *B*-encoding). Second, distributed codes in the property approach are all non-local,
693   because joint *B*-computations across multiple **r**-locations must be completed at a single point in **r**-
694   space, before the result of that computation generates some $<s_{ab}>(\rho_j)$ (as the *property-of* the
695   completed *B*-computation value).

### 7.    Locality requires that point-like states encode components-of-consciousness

697   Developments in Section 6 informally establish that locality requires non-distributed codes for
698   consciousness (because [6.2.1.1] established that sub-domain approaches can only be local if codes
699   are non-distributed, and [6.2.2.2] established that all distributed codes in the property approach are
700   non-local). Although ρ-formalism [5.4] is foundational to the somewhat intuitive approach employed
701   in Section 6, the full *D* symbolism for theories-of-consciousness [5.6] was not employed there. Nor
702   were the demonstrations of Section 6 premised on a *formal* definition of local causality. The present
703   Section remedies these deficits by providing a formal analysis of locality in physically-orthodox
704   theories. Formal proofs are not entirely independent of the intuitive developments of Section 6,
705   however, because they rely on the sub-domain/property bifurcation established there [6.1], as well as
706   on the definition of **r**/ρ-locality in [6.2.2.1]. Nevertheless, formal developments in the present
707   Section supply both an additional degree of *rigor*, and an extension in the *generality* of analyses, to
708   consider parallel and temporal codes deliberately neglected in Section 6. This Section also addresses
709   the key question of why physical orthodoxy excludes distributed *B*-coding of consciousness but does
710   *not* exclude distributed *A*-codes for neural computations serving behavior.

711

### 7.1.    Locality requires that point-like states encode components-of-consciousness

### 7.1.1.  Formal characterization of local causality

Let $P_1$ be a symbol for the measured value of a physical property at a location with coordinates $\chi_1$, and $P_n$ the analogous symbol for some other property at $\chi_n$. Then "locality of causation" means that a change $\Delta P_1$ at $\chi_1$ can only cause a change $\Delta P_n$ at $\chi_n$ if $\chi_1$ and $\chi_n$ are two labels for a single physical location, or if there is a continuous path $\chi_2, \ldots, \chi_{n-1}$ from $\chi_1$ to $\chi_n$ such that $\Delta P_1$ results in a physical change $\Delta P_2$ at $\chi_2$, $\Delta P_2$ results in $\Delta P_3$ at $\chi_3$, and so on, until $\Delta P_{n-1}$ results in the change $\Delta P_n$ at $\chi_n$.

For present purposes, a change in neural state $\Delta A(_{A}\mathbf{r}_1)$ at some location $_{A}\mathbf{r}_1$ can certainly have distant effects, say on muscle activity (*e.g.* after intermediate, behaviorally-relevant, neural computations), if there is a path of physical connectivity from $_{A}\mathbf{r}_1$ to relevant muscle sites. Any such spatially-distant effect is then implemented via a chain of local causes and effects. Consider, in contrast, the effect of $\Delta B(_{B}\mathbf{r}_1)$ on the component of consciousness $<s>$ at $\rho_1$, say. $\Delta B(_{B}\mathbf{r}_1)$ can only engender a change $\Delta<s>$ at $\rho_1$ that is *consistent with local causation* if the coordinates $_{B}\mathbf{r}_1$ and $\rho_1$ label the same physical location, or if there is a connected physical path from $_{B}\mathbf{r}_1$ to $\rho_1$ that also contains causally efficacious physical structures, transmitting $\Delta B(_{B}\mathbf{r}_1)$ to $\Delta<s>$ via a chain of local changes in physical state.

The analysis of local causality in the generation of conscious experience therefore requires some knowledge of the relative location, ontology, and topology of locations labelled by $\mathbf{r}$ coordinates and those labelled by $\rho$ coordinates. Although, as discussed in [5.4], considerable uncertainty currently exists regarding these matters in the *general* case, comprehensive analysis in Section 6 of possibilities in *physically-orthodox* settings will suffice as a firm foundation for the formal analyses in [7.2-7.3].

734

### 7.1.2.  Locality in classical and quantum physical theories

Spatiotemporally-local causality is built into modern physical theory, because almost all phenomena amenable to an orthodox explanation must be reducible to the existence and local interactions of elementary particles and spacetime curvature. Although quantum phenomena are sometimes advanced as a broad exception to local causality, for present purposes, it suffices to note that quantum non-locality (if it exists) is not of an arbitrary kind that can be harnessed to any purpose. Moreover, because the central present focus is to characterize a biophysical test for physically-orthodox theories-of-consciousness that are (as noted in the Introduction) all purely *classical* in nature, further discussion of quantum non-locality will be given as a discussion point [8.1.3], rather than as a feature of central analysis.

745

### 7.2.    Spatial locality

### 7.2.1.  Spatial locality excludes distributed *B*-codes for a single feature type

---

Nicholas M. Rosseinsky

748  The intuitive discussion of consequences from spatial locality for a physically orthodox theory-of-
749  consciousness given in [6.2] can be stated in more rigorous formal terms, by employing the
750  $D_{abj}[\{B(_B\mathbf{r}_i)\}]$ symbolism of Eqs. 14-16. Specifically, consider a single modality (*e.g.* vision) and a
751  distributed code in which activity at multiple locations encodes a single feature type at a single $\rho$
752  location. For example, let activity at three separate locations $_B\mathbf{r}_1$, $_B\mathbf{r}_2$, and $_B\mathbf{r}_3$ encode the conscious
753  experience of color, $<s_{col,b}>$ say, at $\rho_1$. (In this notation, *b* indexes the various different colors that can
754  be experienced at $\rho_1$.) Following Eq. 14, writing $D_{col,b,1}$ for the *B*-classifier relevant to color at $\rho_1$, it
755  follows that

756  $$D_{col,b,1}[\{B(_B\mathbf{r}_1), B(_B\mathbf{r}_2), B(_B\mathbf{r}_3)\}] = 1 \Leftrightarrow \exists_P <s_{col,b}>(\rho_1) \tag{17}.$$

757  Locality problems are immediately apparent in Eq. 17, because the LHS requires integration of *B*-
758  information at multiple locations that are not connected by an $<s_{col,b}>(\rho_1)$-relevant causal physical
759  structure [7.1.1].

760

### 7.2.1.1.     Case 1: sub-domain scheme

762  Consider the case in which $\rho$-coordinates are labels for a sub-domain of orthodox physical three-
763  dimensional space [6.1,6.2.1], so that points labeled by $_B\mathbf{r}_1$, $_B\mathbf{r}_2$, $_B\mathbf{r}_3$, and $\rho_1$ can all be viewed as lying
764  in a single, physically-connected, space. Consider further the *degenerate* case in which $\rho_1$ labels a
765  point that is already labeled by one of the $\{_B\mathbf{r}_1, _B\mathbf{r}_2, _B\mathbf{r}_3\}$, $_B\mathbf{r}_1$ say (without loss of generality), *i.e.*
766  where the point in orthodox three-dimensional space used to generate $<s_{col,b}>(\rho_1)$ is *the same as* one
767  of the points used for *B*-measurement of $D_{col,b,1}$ conditions. (The same result follows trivially for the
768  *non-degenerate* or *general* sub-domain case, where $\{_B\mathbf{r}_1, _B\mathbf{r}_2, _B\mathbf{r}_3, \rho_1\}$ all label different points.) In the
769  degenerate case, information must be propagated from $_B\mathbf{r}_2$ and $_B\mathbf{r}_3$ to $\rho_1 = _B\mathbf{r}_1$ in order to assess
770  $D_{col,b,1}$ values, and the lack of intermediate physical structure for this propagation means that the
771  generation of $<s_{col,b}>(\rho_1)$ is *not* local, according to [7.1.1].

772      Neural connectivity between $\{_B\mathbf{r}_1, _B\mathbf{r}_2, _B\mathbf{r}_3\}$ and $\rho_1$ cannot be the intermediate mechanism for
773  propagating $D_{col,b,1}$ information. For example, consider the general case, and assume that *B*-measures
774  at $\{_B\mathbf{r}_1, _B\mathbf{r}_2, _B\mathbf{r}_3\}$ were propagated to $\rho_1$ by orthodox physical and neural means (*i.e.* by a chain of
775  physically-orthodox local causes and effects). In this case, the encoding of $<s_{col,b}>(\rho_1)$ is *not*
776  *distributed*, contradicting the assumption in [7.2.1], because under these circumstances the correct
777  theory of consciousness is

778  $$D_{col,b,1}[B(_B\mathbf{r}_4)] = 1 \Leftrightarrow \exists_P <s_{col,b}>(\rho_1) \tag{18}$$

779  replacing Eq. 17, where $_B\mathbf{r}_4$ is the **r**-coordinate label for the point also labeled by $\rho_1$. (Eq. 18 must
780  apply under these circumstances, because *B* was defined [5.5.3] as the final informational basis of
781  consciousness in brain dynamics.)

782

### 7.2.1.2.     Case 2: property scheme

784  In a physically orthodox setting, the only alternative to the sub-domain interpretation [6.1,6.2.1] of $\rho$-

785   space is the property scheme [6.1,6.2.2]. In this conception, points labelled by $\{_B\mathbf{r}_1, {}_B\mathbf{r}_2, {}_B\mathbf{r}_3\}$ exist in
786   a common space (the physically-orthodox three-dimensional space), whereas $\rho$-coordinates label
787   points in a separate (physically unconnected) space that is itself the property of brain dynamics. Even
788   if $\mathbf{r}/\rho$-locality is admitted [6.2.2.1], $\{_B\mathbf{r}_1, {}_B\mathbf{r}_2, {}_B\mathbf{r}_3\}$-located information must first be integrated to
789   some point $_B\mathbf{r}_4$, via a definitively $\mathbf{r}$-non-local mechanism, before acting at $\rho_1$ via a putatively local
790   (not-non-local) $_B\mathbf{r}_4$-$\rho_1$ connection.

### 7.2.2.  Locality does not exclude distributed *A*-coding for behavioral processes

793   Taken together, [7.2.1.1-7.2.1.2] imply that the informational basis for any particular component-of-
794   consciousness $<s_{a,b}>$ must be spatially point-like, *i.e.* depend only on a *single* $_B\mathbf{r}$ coordinate. This
795   means that the encoding for consciousness of a particular feature-instance cannot be distributed, and
796   must be contained in a single neuron.

797   However, these results do *not* contradict *e.g.* the well-established roles for distributed codes in
798   computations that serve the generation of behavior (as opposed to the generation of consciousness).
799   Provided that "wiring" (physical connectivity of an appropriate kind [7.1.1]) links all locations
800   involved in a distributed code with causally-later information processing machinery, local causality is
801   not violated, in these contexts. It is precisely the absence of such wiring in the case of consciousness
802   that leads to the conflict between distributed codes and local causality.

### 7.2.3.  Locality does not exclude parallel *B*-processing of distinct feature types

805   It is perhaps tempting to conclude from [7.2.1.1-7.2.1.2] that *all* feature types at a single $\rho$-location
806   must be encoded by a single neuron, *i.e.* must depend only on a *single* $_B\mathbf{r}$ coordinate. This would
807   mean that the generation of conscious experience at a given location $\rho_i$ would have to obey

808   $$D_W[B(_B\mathbf{r}_i)] = 1 \Leftrightarrow \{<s_{a,b}>(\rho_i), (a,b) \in W\} \tag{19}$$

809   where $W$ is a set of $(a,b)$ pairs defining the totality of features at $\rho_i$, and the member constitution of $W$
810   is determined solely by the *B*-measure at a single brain location $_B\mathbf{r}_i$. (Recall that a one-to-one
811   correspondence between the $\{\rho_i\}$ and the $\{_X\mathbf{r}_i\}$ was declared in [5.4]. Eq. 19 defines a new one-to-
812   one relationship between the $\{\rho_i\}$ and the $\{_B\mathbf{r}_i\}$.)

813   However, reasoning in [7.2.1.1-7.2.1.2] that considers a single feature type could only be extended
814   to all feature types simultaneously if it were known that a single physical computation occurs for
815   each $\rho$-location to determine all features at that location. Consider the alternative, namely that there is
816   a set of computations, one for each feature type. Following [7.2.1.1-7.2.1.2], let $_B\mathbf{r}_{col,1}$ be the single
817   *B*-location encoding color at $\rho_1$, and let $_B\mathbf{r}_{edge,1}$ be the single *B*-location encoding edge-orientation at
818   $\rho_1$. For simplicity, consider the property scheme [6.1,6.2.2]. Here, it is possible for the equations

819   $$D_{col,b,1}[B(_B\mathbf{r}_{col,1})] = 1 \Leftrightarrow \exists_P <s_{col,b}>(\rho_1) \tag{20}$$

820   and

---

Nicholas M. Rosseinsky

821 $$D_{\text{edge},c,1}[B(_{\text{B}}\mathbf{r}_{\text{edge},1})] = 1 \Leftrightarrow \exists_{\text{P}} <s_{edge,c}>(\boldsymbol{\rho}_1) \qquad (21)$$

822   to pertain, without contradicting locality, if $<s_{col,b}>(\boldsymbol{\rho}_1)$ is interpreted as a property of $B(_{\text{B}}\mathbf{r}_{\text{col},1})$, and
823   $<s_{edge,c}>(\boldsymbol{\rho}_1)$ as a property of $B(_{\text{B}}\mathbf{r}_{\text{edge},1})$. (This interpretation relies on the definitional separation [3.1]
824   of $s_{col,b}$ and $s_{edge,c}$, and the consequent fact that there are two separate $D$ equations, *i.e.* Eqs. 20-21 . It
825   cannot be applied to give a local interpretation of a single-feature distributed code such as Eq. 17,
826   because self-evidently this equation inherently contains multi-$\mathbf{r}$ references.)

827   In summary, if local causality applies, the $a$-th feature type must be encoded by a non-distributed
828   code that defines the $x$-instance $<s_{a,x}>$ for a given $\boldsymbol{\rho}_j$ by the $B$-value at a *single* $\mathbf{r}$ location, $_{\text{B}}\mathbf{r}_{a,j}$ say. But
829   various different feature types $s_a$, $s_b$, $s_c$, and so on, can have their instances $<s_{a,x}> <s_{b,y}> <s_{c,z}>$ at $\boldsymbol{\rho}_j$
830   defined by $B$-values at a *different* $\mathbf{r}$ locations $_{\text{B}}\mathbf{r}_{a,j}$, $_{\text{B}}\mathbf{r}_{b,j}$, $_{\text{B}}\mathbf{r}_{c,j}$ *etc.* in a manner consistent with local
831   causality, *e.g.* under certain interpretations of the property conception of $\boldsymbol{\rho}$-space.

832

### 7.2.4. Locality does not exclude topographic $B$-processing for a single feature type

834   Emphatically, results in [7.2.1.1-7.2.1.2] do not contradict the feasibility of parallel (as opposed to
835   distributed) processing of a *single* feature in physically-orthodox theories-of-consciousness. Thus, for
836   example, $<s_{1,b(1)}>(\boldsymbol{\rho}_1)$, $<s_{1,b(2)}>(\boldsymbol{\rho}_2)$, $<s_{1,b(3)}>(\boldsymbol{\rho}_3)$ and so on can be encoded as the result of parallel
837   computations that generate $B$-values at $_{\text{B}}\mathbf{r}_{1,1}$, $_{\text{B}}\mathbf{r}_{1,2,}$ $_{\text{B}}\mathbf{r}_{1,3}$ *etc.*, with

838   $$D_{1,b(j),j}[B(_{\text{B}}\mathbf{r}_{1,j})] = 1 \Leftrightarrow \exists_{\text{P}} <s_{1,\,b(j)}>(\boldsymbol{\rho}_j) \qquad (22).$$

839   Eq. 22 can, for example, describe a (non-distributed) *topographic* encoding. Arguments of [7.2.1.1-
840   7.2.1.2] exclude only (and importantly) the appearance of multiple $\mathbf{r}$-locations in the argument of $D$
841   for a given $<s_{a,b}>(\boldsymbol{\rho}_j)$, *i.e.* terms of the form $D_{abj}[\{B(_{\text{B}}\mathbf{r}_i)\}]$ (recalling from [3.2] that $\{\ldots\}$ notation in
842   the present work is reserved for multi-member sets or collections).

843

### 7.2.5. Spatially-extended phenomena in orthodox physics

845   One potential confusion concerning results so far must be cleared up. Specifically, it might appear
846   that orthodox physical theory contains spatially-extended properties without intermediate causal
847   structure (**Figure 12**). For example, temperature is a physically orthodox property that appears to
848   pertain to molecules that are spatially distributed, so it might appear that consciousness could
849   similarly be a property of a spatially-distributed code, without violating physical orthodoxy.
850   However, careful analysis of temperature shows that it is a *measurement* involving a *spatially-*
851   *extended measuring apparatus* (thermometer). Apparatus of this kind constitutes precisely the kind
852   of "intermediate causal structure" [7.2.6] that is absent in a physically-orthodox theory-of-
853   consciousness. Thus analogies between temperature and consciousness are not valid, with respect to
854   spatial extension.

855

### 7.2.6. Roles of "3-D space" and "conventional particle spectrum" assumptions

857   In [5.2], physically-orthodox theories-of-consciousness were specifically restricted to three-

858   dimensional space and a standard particle spectrum. These restrictions implicitly support the
859   arguments of [7.2.1.1-7.2.1.2] by ruling out the *conceptual* possibility of a physical structure beyond
860   the brain that could (even in principle) link encoding sites $\{_B\mathbf{r}_i\}$ to $\{\rho_j\}$ locations (**Figure 12D**).

861       If space were *e.g.* six- rather than three-dimensional, one could conceive of {quark, electron,
862   photon}-constituted physical structures linking encoding sites $\{_B\mathbf{r}_i\}$ in the brain with $\{\rho_j\}$ locations in
863   a three-dimensional subspace orthogonal to that containing the brain. (Extra dimensionality is
864   required if these structures are constituted of orthodox particles, because such structures would be
865   directly observed if they existed in the three-dimensional subspace containing the brain.)

866       If non-standard particles could bind together to constitute physical structures linking encoding
867   sites $\{_B\mathbf{r}_i\}$ in the brain with $\{\rho_j\}$ locations, and these structures had no interaction with orthodox
868   matter other than at $\{_B\mathbf{r}_i\}$ sites, then extra-dimensionality would not be required to house information
869   transmission mechanisms.

870       Thus, excluding extra-dimensionality and non-orthodox particles (as is conventional, in physically
871   orthodox theories) rules out certain hypothetical possibilities that could contradict the "no
872   intermediate causal structure" claims of [7.2.1.1-7.2.1.2]. Conversely, as will be discussed in [8.1.2],
873   explicitly introducing extra-dimensionality or non-orthodox particles can allow distributed codes to
874   be consistent with local causation, by providing the means for intermediate causal structures.

875

876   **7.3.    Temporal locality**

877   **7.3.1.  Temporal locality excludes rate and temporal codes**

878   A rate code (Adrian, 1928) necessarily requires the employment of a temporally extended interval
879   over which the rate or frequency of firing is established. More complex temporal codes (Reinagel and
880   Reid, 2000; Mehta *et al.*, 2002; Panzeri *et al.*, 2010) also require a finite (non-zero) temporal interval
881   for decoding. For present purposes, extended temporal decoding intervals are analogous to the spatial
882   distribution of *B*-dynamics, in the sense that they create an informational basis for consciousness that
883   is *not* spatiotemporally point-like.

884       Consider a single neuron rate-code for $<s_{a,b}>(\rho_j)$ in which *B* activity must be assessed over a
885   temporal interval $\{t, \ldots, t + \varepsilon\}$ (where $t$ is a conventional time coordinate, and $\{t, \ldots, t + \varepsilon\}$ denotes
886   the continuum of points from $t$ to $t + \varepsilon$, including the end points). Writing $B(\mathbf{r},t)$ for the time-$t$ value
887   of *B* at $\mathbf{r}$, the theory-of-consciousness for such a code must be of the form

888         $$D_{abj}[B(_B\mathbf{r}_i, \{t,\ldots, t + \varepsilon\})] = 1 \Leftrightarrow <s_{a,b}>(\rho_j, \{t + \varepsilon,\ldots, t + \varepsilon + \Delta\})) \qquad (23)$$

889   where $<s_{a,b}>(\rho_j, t)$ denotes the time-$t$ existence of $<s_{a,b}>(\rho_j)$, and the time interval on the RHS
890   describes the existence of $<s_{a,b}>(\rho_j)$ for a period of length $\Delta$ beginning at $t + \varepsilon$. [For example, $\Delta$ is
891   order $10^{-1}$ s, if conscious experience is refreshed every 100 ms or so (VanRullen and Koch, 2003).]
892   The absence of intermediate causal structure to act as a kind of memory, and thus effect temporally
893   extended processing, means that a locally-causal *D* cannot act on a temporally-extended set of *B*
894   values in a physically orthodox setting, thus excluding Eq. 23-based theories and a wide range of
895   temporal codes, including rate codes.

Nicholas M. Rosseinsky

896    Exclusion of spatially-distributed and temporally-extended codes means that a particular feature
897    $<s_{a,b}>$ at a given coordinate location $\rho_j$ must be encoded by the instantaneous value of $B$ at a single
898    point in space and time, for local causality, so that

899         $$D_{abj}[B(_B\mathbf{r}_i, t)] = 1 \Leftrightarrow <s_{a,b}>(\rho_j, t) \qquad (24).$$

900    This does not mean that $<s_{a,b}>(\rho_j)$ itself cannot be temporally extended, but only that such temporal
901    extension must be generated by temporally-extended satisfaction of the condition on the LHS of Eq.
902    24. That is, in order to generate $<s_{a,b}>(\rho_j, \{t_0 ,…, t_0 + \Delta \})$, we must have $D_{abj}[B(_B\mathbf{r}_i, t)] = 1$ for each $t$
903    in $\{t_0 ,…, t_0 + \Delta \}$. (Note the difference between this condition and the LHS of Eq. 23. In Eq. 24, the
904    satisfaction of $D_{abj}[B(_B\mathbf{r}_i, t)] = 1$ for each $t$ has the temporally-local effect that $<s_{a,b}>(\rho_j, t)$ exists. In
905    contrast, in Eq. 23 the action of $D$ requires first the temporally *non*-local accumulation of $B$-values
906    throughout the interval $\{t,…, t + \varepsilon\}$.)

907    Eq. 24 challenges a role for spike-based codes as the informational basis of consciousness in a
908    locally-causal theory, because a neuron can only be in the spiking state for a few milliseconds at a
909    time. (That is, even during a temporally-extended high frequency burst, spikes are interspersed with
910    non-spike intervals). Thus, temporally local generation of a component-of-consciousness that itself
911    persists continuously for 100 ms or more cannot be based on spike states, because spike states only
912    persist for a few milliseconds. (At best, if firing were sustained for 100 ms, temporally-local
913    generation must mean that the component-of-consciousness "flickers" in and out of existence in
914    precise synchronization with the spiking/non-spiking state of the neuron.)

915    Despite the inability of spike-based codes to produce temporally continuous and persistent
916    components of consciousness, other mechanisms are possible. Consider, for example, a neuron that
917    encodes a feature $s_{a,b}(_X\mathbf{r}_j)$ by spiking activity of some kind, which in turn puts a biomolecule $M_{abj}$
918    [5.5.2] into some energy state $E[M_{abj}](t_0)$ at time $t_0$. (In this notation, $E[O]$ is the energy of a physical
919    object labelled by the symbol $O$, and $E[O](t)$ is the energy of $O$ at time $t$.) Assume that the classifier
920    $D_{abj}$ for $<s_{a,b}>(\rho_j)$ can be written as

921         $$E[M_{abj}](t) > E_0 \Leftrightarrow <s_{a,b}>(\rho_j, t) \qquad (25),$$

922    *i.e.* that $<s_{a,b}>(\rho_j)$ exists when the energy of $M_{abj}$ is above some threshold $E_0$. Then, if the burst-
923    induced energy $E[M_{abj}](t_0)$ is greater than $E_0$, and if a time interval of duration $\Delta$ elapses before
924    dissipative mechanisms result in $E[M_{abj}]$ falling below $E_0$, it follows from Eq. 25 that $<s_{a,b}>(\rho_j)$ exists
925    continuously for the period $<s_{a,b}>(\rho_j, \{t_0 ,…, t_0 + \Delta \})$.

926

927    ### 7.3.2. Temporally-extended phenomena in orthodox physics

928    Just as certain orthodox physical properties might appear to be spatially extended [7.2.5], the
929    appearance of temporal extension in physically-orthodox settings might seem to invalidate the claim
930    that physically-orthodox codes must be temporally point-like. Consider, for example, the simple
931    temporal persistence of a rigid body (in the absence of conditions that would cause it to move or to
932    fragment). Clearly, although temporal extension is real, it results from a succession of temporally-
933    local persistences of component elementary-particles. In contrast, the conversion of a temporally-
934    extended code (*e.g.* Eq. 23) involved in the generation of some component-of-consciousness requires

935  non-temporally-local *memory* of physical features, meaning that the analogy between temporal
936  extension of the rigid body and temporally-extended computation breaks down.

937

## 8.    Discussion

### 8.1.    Experimental tests and implications from future experimental results

#### 8.1.1.    Biophysical structures must convert distributed to local codes

941  Results in [7.2] and [7.3] mean that physically-orthodox theories-of-consciousness must propose
942  spatiotemporally point-like encodings of contributions-to-consciousness, if these theories are to be
943  consistent with locality of physical causation. This implication contrasts with the widespread
944  prevalence of spatially-distributed (Georgopoulos *et al*., 1986; Buonomano and Merzenich, 1995;
945  Tsodyks *et al*., 1996; Pillow *et al*., 2008; Solomon and Lennie, 2007; Quiroga and Panzeri, 2009) and
946  temporally extended (Adrian, 1928; Reinagel and Reid, 2000; Mehta *et al*., 2002; Panzeri *et al*.,
947  2010) encoding schemes in current accounts of neural information processing that serves behavior.
948  Point-like encoding of consciousness and distributed/extended behavioral codes can only be
949  reconciled if biophysical mechanisms exist in the brain to convert non-point-like codes to
950  spatiotemporally local equivalents, for each encoding that is used as the informational basis of
951  consciousness. These mechanisms must be empirically observable, thus creating a strong
952  experimental constraint for locally-causal generation of contributions-to-consciousness based on
953  brain dynamical states. Potential methods for verification or falsification of these predictions include
954  detailed neuroanatomical investigations of dynamically-relevant brain structure (Alivisatos *et al*.,
955  2013) that might uncover previously unknown physical connectivity, and detailed simulation models
956  (Markram, 2006) that might uncover previously unappreciated biophysical mechanisms for point-like
957  integration of non-local codes.

958

#### 8.1.2.    Consequences from absence of distributed-to-local conversion

960  If biophysical mechanisms for integration of distributed codes can be definitively excluded by future
961  anatomical and computational investigations, physically orthodox theories cannot explain
962  consciousness. Accordingly, any scientific explanation of consciousness must relax one or more the
963  defining features of physical orthodoxy [5.2]. Most directly, the requirement for locality in physical
964  causation can be dropped. However, this leads to another problem: ensuing non-locality is of a very
965  particular kind, with a high degree of order – certain sets of locations in the brain are connected
966  together computationally, and then connected to specific $\rho$-coordinate locations. Without physical
967  structure to effect these connections, the details of connections appear to be inexplicable, thus
968  rendering consciousness itself beyond scientific explanation.

969      As noted in [7.2.6], relaxation of the three-spatial-dimensions or conventional-particle-spectrum
970  assumptions (rather than the relaxation of the locality assumption) can lead to theories that both
971  explain the complex order of connections and preserve locality of physical causation. These theories
972  can hypothesize connective physical structures that are not observable via conventional means, and
973  preserve the possibility of a scientific explanation of consciousness. Thus, if biophysical mechanisms
974  for reduction to local codes are excluded, results here indicate that higher-dimensional or non-

975    orthodox particle spectrum theories-of-consciousness should be prioritized, in the first instance.

976

### 8.1.3.  "Quantum consciousness" and distributed-to-local conversion

978    Informal discussions of locality issues elsewhere have resulted in observations such as "any physical
979    process responsible for consciousness would have to be something with an essentially global [non-
980    local] character" (Penrose, 2000), and a related inference that consciousness must therefore explicitly
981    depend on quantum mechanisms. More precise, symbol-based, formal considerations in the present
982    paper provide a number of grounds for challenging the alleged *necessity* for an explicitly quantum
983    theory-of-consciousness based on hypothetical non-locality in brain action on a number of grounds.
984    First, it is possible that biophysical mechanisms exist to convert non-local to local codes (and will be
985    observed) [8.1.1]. Second, if these mechanisms are excluded, at least two classical proposals
986    (involving physical structure either in extra dimensions or comprised of non-orthodox particles) exist
987    [8.1.2] as alternatives to "quantum consciousness". Finally, even if quantum non-locality is involved
988    *e.g.* in the integration of spatially-distributed brain-dynamical information, further mechanisms
989    explaining the precise order of connectivity [8.1.2] amongst $\{_B\mathbf{r}_i\}$ locations and between $\{_B\mathbf{r}_i\}$ and
990    $\{\rho_j\}$ locations must be proposed: these do not follow merely from the hypothesis of a quantum
991    mechanism, and are signally lacking in current quantum theories-of-consciousness [*e.g.* (Penrose,
992    1989; Stapp, 1993; Beck and Eccles, 2003)].

993

### 8.2.    Robustness

995    The exclusion of distributed codes for consciousness is hard to avoid, precisely because the general
996    approach creates few assumptions whose invalidity might falsify reasoning. Note that an objection
997    such as "neural coding of consciousness might not depend on distributed codes" does not refute
998    results: it is simply another way of stating that biophysical mechanisms must exist to reduce
999    distributed or extended codes to spatiotemporally point-like states.

1000       If psychophysical parallelism is excluded [2.3], the only apparent way to avoid the reduction of
1001    distributed codes for a locally-causal theory is to deny that conscious experience exists *anywhere*.
1002    This is apparently the route taken by Dennett (Dennett, 1991), although this position is, at best,
1003    difficult to understand, because it seems to mean that conscious experience itself does not exist. In
1004    the present paper, non-existence is excluded by the presumption that each reader can verify for
1005    themselves the meaning of the terminology "exteroceptive, visual and auditory, conscious
1006    experience" [2.1,2.2].

1007

### 8.3.    Outlook

1009    Central results here help to shape the outlook for the scientific understanding of human
1010    consciousness. Because neither spatially-distributed nor temporally-extended codes can be
1011    physically-orthodox informational bases for consciousness, point-like encodings and associated
1012    biophysical machinery for their generation must be observed if the currently mainstream physically-
1013    orthodox presumption is to be sustained, thus establishing for the first time experimental tests for the

1014    physically-orthodox hypothesis. Emphatically, tests here do not arbitrate between alternative "hows"
1015    for consciousness (such as phase synchrony, or re-entrant processing); instead, they address the more
1016    fundamental issue of physical orthodoxy that has until now been amenable only to philosophical
1017    (McGinn, 1989; Dennett, 1991; Chalmers, 1996) or even mathematical (Penrose, 1989) analysis.

1018    Experimental outcomes can provide more than (currently absent) empirical support for the physical
1019    orthodoxy presumption, however. Observations supporting physical orthodoxy should inform future
1020    theoretical and experimental work, by identifying brain locations and physical properties critically
1021    involved in the generation of consciousness. Of course, observations rejecting physical orthodoxy
1022    would advance the field by ruling out a large class of theories, but theoretical analyses here *also*
1023    identify a small group of basic theories that would then become prominent. Notably for the rational
1024    conduct of a science of consciousness, even after a putative rejection of physical orthodoxy, the
1025    content and multiplicity of theories then remaining rejects allegedly conclusive assertions that
1026    consciousness "must" be quantum (Penrose, 1989; Stapp, 1993; Beck and Eccles, 2003), non-
1027    physical (Chalmers, 1996), or beyond human understanding (McGinn, 1989), because remaining
1028    candidates are all physical (albeit of a non-orthodox kind) and contain no definitively quantum
1029    features. Hence, in addition to the characterization of novel and important experiments, developments
1030    here can contribute a new dialectical order to the field as whole, helping to instil a clarity that can
1031    only be positive for the science of consciousness.

1032    Claims that present methods can provide basic clarity are founded in part on precision due to new
1033    formal symbolism for describing phenomena and expressing theoretical alternatives. For example,
1034    symbolism has been employed in this paper to treat in precise terms: the delineation between binding
1035    problems in behavioral computations (described in $A$ and $\mathbf{r}$ notation) and in conscious experience ($B$,
1036    $\rho$ notation); the relationship between orthodox and conscious-experiential spaces (labelled by $\mathbf{r}$ and $\rho$
1037    symbols respectively); and, both spatial and temporal locality. [(Rosseinsky, 2014b, 2014a) apply
1038    symbolism developed here to give a fuller treatment of binding problems.] Notably, a formal
1039    approach resolves previously imprecise verbal appeals to emergent properties of complex systems, by
1040    revealing the role that intermediate physical structure must play if these appeals are to explain
1041    spatially- and temporally-extended phenomena. Precision of this kind is essential if consciousness is
1042    to be explained scientifically.

1043    Relatedly, the present approach supports methodological clarity by emphasizing: explicit definitions
1044    of phenomenological scope; definitions of physical orthodoxy and its requirements; explicit
1045    statement of theories considered and excluded; and explicit demarcation between (large classes of)
1046    theories whose comparative status is to be assessed by experiments. Proper treatment of these
1047    methodological basics assuages concerns that the study of consciousness is not scientific (Wilkes,
1048    1984, 1988; Hawking, 2000). (Related concerns stemming from the role of subjective report are
1049    minimized here by experiments that are entirely *objective* in nature.)

1050    The primary question here has been the conditions under which parallel and distributed codes can
1051    generate components-of-consciousness in a manner obeying the local causality requirement of
1052    physical orthodoxy. Central results exclude a physically-orthodox encoding-of-consciousness role for
1053    spatially-distributed and temporally-extended codes under *any* conditions, but *unconditionally* permit
1054    such a role for parallel codes. However, physical orthodoxy places more requirements on theories
1055    than just local causality. Subsequent papers in the present series (Rosseinsky, 2014b, 2014a) employ
1056    the basic approach developed here to investigate further restrictions on parallel codes imposed by
1057    physical orthodoxy, extending the range of experimental tests for physically-orthodox theories-of-

Nicholas M. Rosseinsky

1058　consciousness, and providing further demonstration of the positive contributions available from the
1059　approach initiated in this paper.

1060

1061    **Table 1. Summary of formal symbols.**

| SYMBOL | DEFINITION | SECTION |
|---|---|---|
| $s_{ab}$ | $b$-th instance of $a$-th feature-type (stimulus) in the external environment. | 3.1 |
| $_X\mathbf{r}_j$ | $j$-th location in the external environment for sampling of sensory information | 3.3 |
| $<s_{ab}>$ | Contribution to conscious experience generated by $s_{ab}$ | 5.3 |
| $\rho_j$ | $j$-th location in conscious experience; information sampled at $_X\mathbf{r}_j$ leads to a contribution-to-consciousness at $\rho_j$ | 5.4 |
| $A$ | Measure of brain activity for behavioral encoding | 4.2 |
| $B$ | Measure of brain activity for encoding and final brain-dynamical cause of consciousness | 5.5.3 |
| $_A\mathbf{r}_i$ | $i$-th brain location relevant to $A$ measurement | 4.4 |
| $_B\mathbf{r}_i$ | $i$-th brain location relevant to $B$ measurement | 5.5.3 |
| $C_{abj}$ | Classifier function on $A$-states for behavioral encoding | 4.5 |
| $D_{abj}$ | Classifier function on $B$-states for encoding and final brain-dynamical cause of consciousness | 5.6 |
| $\exists_P$ | Denotes existence of a physical object, property, or phenomenon | 4.5 |
| $\nexists_P$ | Denotes absence of a physical object, property, or phenomenon | 4.5 |
| $\{ …\}$ | Set or collection of … | 3.2 |

1062

1063

## 9.     References

Abbott, L., and Sejnowski, T. J. (1999). *Neural codes and distributed representations: foundations of neural computation*. Cambridge, MA: MIT Press.

Abeles, M. (2011). Cell assemblies. *Scholarpedia* 6, 1505. doi:10.4249/scholarpedia.1505.

Adrian, E. D. (1928). *The basis of sensation: the action of the sense organs*. New York: W. W. Norton & Co.

Alivisatos, A. P., Chun, M., Church, G. M., Deisseroth, K., Donoghue, J. P., Greenspan, R. J., McEuen, P. L., Roukes, M., Sejnowski, T. J., and Weiss, P. S. (2013). The brain activity map. *Science* 339, 1284–1285.

Anderson, P. W. (1972). More is different. *Science* 177, 393–396.

Baars, B. J. (1988). *A Cognitive Theory of Consciousness*. Cambridge: Cambridge University Press.

Baars, B. J., and Edelman, D. B. (2012). Consciousness, biology and quantum hypotheses. *Phys. Life Rev.* 9, 285–294.

Beck, F., and Eccles, J. C. (2003). Quantum processes in the brain: A scientific basis of consciousness. *Adv. Conscious. Res.* 49, 141–166.

Bell, A. J., and Sejnowski, T. J. (1996). The" independent components" of natural scenes are edge filters. *Vision Res.* 37, 3327–3338.

Buonomano, D. V., and Merzenich, M. M. (1995). Temporal information transformed into a spatial code by a neural network with realistic properties. *Science*, 1028–1030.

Buzsáki, G. (2006). *Rhythms of the brain*. Oxford; New York: Oxford University Press.

Chalmers, D. J. (1996). *The Conscious Mind: In Search of a Fundamental Theory*. New York: Oxford University Press.

Chalmers, D. J. (2008). "The varieties of emergence," in *The Re-emergence of emergence: The Emergentist Hypothesis from Science to Religion*, eds. P. Davies and P. Clayton (Oxford: Oxford University Press), 244–256.

Churchland, P. M. (1985). Reduction, qualia, and the direct introspection of brain states. *J. Philos.* 82, 8–28.

Churchland, P. S. (2005). A neurophilosophical slant on consciousness research. *Prog. Brain Res.* 149, 285.

Churchland, P. S., and Sejnowski, T. J. (1992). *The Computational Brain*. Cambridge, MA: MIT press.

Crick, F., and Koch, C. (1990). Towards a neurobiological theory of consciousness. *Semin. Neurosci.*

1096        2, 263–275.

1097    Decharms, R. C., and Zador, A. (2000). Neural representation and the cortical code. *Annu. Rev.*
1098        *Neurosci.* 23, 613–647.

1099    Dehaene, S. (2014). *Consciousness and the brain: deciphering how the brain codes our thoughts.*
1100        New York: Viking.

1101    Dennett, D. C. (1991). *Consciousness Explained.* Boston, MA: Little Brown.

1102    Dennett, D. C., and Kinsbourne, M. (1992). Time and the Observer: the Where and When of
1103        Consciousness in the Brain. *Behav. Brain Sci.* 15, 183–247.

1104    Edelman, G. M. (1989). *The remembered present: A biological theory of consciousness.* New York:
1105        Basic Books.

1106    Eggermont, J. J. (1998). Is there a neural code? *Neurosci. Biobehav. Rev.* 22, 355–370.

1107    Engel, A. K., and Singer, W. (2001). Temporal binding and the neural correlates of sensory
1108        awareness. *Trends Cogn. Sci.* 5, 16–25.

1109    Faisal, A. A., Selen, L. P. J., and Wolpert, D. M. (2008). Noise in the nervous system. *Nat. Rev.*
1110        *Neurosci.* 9, 292–303.

1111    Feldman, J. (2013). The neural binding problem(s). *Cogn. Neurodyn.* 7, 1–11.

1112    Field, D. J. (1994). What is the goal of sensory coding? *Neural Comput.* 6, 559–601.

1113    Foley, J. M. (1978). "Primary distance perception," in *Perception* (Berlin: Springer), 181–213.

1114    Foley, J. M., Ribeiro-Filho, N. P., and Da Silva, J. A. (2004). Visual perception of extent and the
1115        geometry of visual space. *Vision Res.* 44, 147–156.

1116    Georgopoulos, A. P., Schwartz, A. B., and Kettner, R. E. (1986). Neuronal population coding of
1117        movement direction. *Science* 233, 1416–1419.

1118    Van Gulick, R. (2001). Reduction, emergence and other recent options on the mind/body problem. A
1119        philosophic overview. *J. Conscious. Stud.* 8, 1–34.

1120    Hameroff, S. (2012). How quantum brain biology can rescue conscious free will. *Front. Integr.*
1121        *Neurosci.* 6. doi:10.3389/fnint.2012.00093.

1122    Hawking, S. W. (2000). "The Objections of an Unashamed Reductionist," in *The Large, the Small*
1123        *and the Human Mind*, ed. M. Longair (Cambridge: Cambridge University Press), 169–172.

1124    Haynes, J.-D., and Rees, G. (2006). Decoding mental states from brain activity in humans. *Nat. Rev.*
1125        *Neurosci.* 7, 523–534.

1126    Heller, J. (1997). On the psychophysics of binocular space perception. *J. Math. Psychol.* 41, 29–43.

Nicholas M. Rosseinsky

1127    Horikawa, T., Tamaki, M., Miyawaki, Y., and Kamitani, Y. (2013). Neural Decoding of Visual
1128        Imagery During Sleep. *Science* 340, 639–642. doi:10.1126/science.1234330.

1129    Judd, D. B., and Wyszecki, G. (1975). *Color in business, science and industry*. Chichester: Wiley.

1130    Kaas, J. H. (1997). Topographic maps are fundamental to sensory processing. *Brain Res. Bull.* 44,
1131        107–112.

1132    Kandel, E. R., Schwartz, J. H., Jessell, T. M., Siegelbaum, S. A., and Hudspeth, A. J. eds. (2012).
1133        *Principles of Neural Science*. New York: McGraw-Hill Professional.

1134    Klimesch, W., Doppelmayr, M., Yonelinas, A., Kroll, N. E. A., Lazzara, M., Röhm, D., and Gruber,
1135        W. (2001). Theta synchronization during episodic retrieval: neural correlates of conscious
1136        awareness. *Cogn. Brain Res.* 12, 33–38.

1137    Di Lollo, V. (2012). The feature-binding problem is an ill-posed problem. *Trends Cogn. Sci.* 16, 317–
1138        321.

1139    Luneberg, R. K. (1944). Mathematical theory of optics (lecture notes). *Brown Univ.*

1140    Mainen, Z. F., and Sejnowski, T. J. (1995). Reliability of spike timing in neocortical neurons. *Science*
1141        268, 1503–1506.

1142    Von der Malsburg, C. (1981). The Correlation Theory of Brain Function. Max Planck Institute for
1143        Biophysical Chemistry, Dept. of Neurobiology.

1144    Markram, H. (2006). The blue brain project. *Nat. Rev. Neurosci.* 7, 153–160.

1145    McFadden, J. (2002). The Conscious Electromagnetic Information (Cemi) Field Theory: The Hard
1146        Problem Made Easy? *J. Conscious. Stud.* 9, 45–60.

1147    McGinn, C. (1989). Can We Solve the Mind–Body Problem? *Mind* 98, 349–366.

1148    Mehta, M. R., Lee, A. K., and Wilson, M. A. (2002). Role of experience and oscillations in
1149        transforming a rate code into a temporal code. *Nature* 417, 741–746.
1150        doi:10.1038/nature00807.

1151    Melloni, L., Molina, C., Pena, M., Torres, D., Singer, W., and Rodriguez, E. (2007). Synchronization
1152        of neural activity across cortical areas correlates with conscious perception. *J. Neurosci.* 27,
1153        2858–2865.

1154    Merker, B. (2007). Consciousness without a cerebral cortex: A challenge for neuroscience and
1155        medicine. *Behav. Brain Sci.* 30, 63–81.

1156    Miller, G. (2005). What Is the Biological Basis of Consciousness? *Science* 309, 79.
1157        doi:10.1126/science.309.5731.79.

1158    Olshausen, B. A. (1996). Emergence of simple-cell receptive field properties by learning a sparse
1159        code for natural images. *Nature* 381, 607–609.

1160 Panzeri, S., Brunel, N., Logothetis, N. K., and Kayser, C. (2010). Sensory neural codes using
1161     multiplexed temporal scales. *Trends Neurosci.* 33, 111.

1162 Penrose, R. (1989). *The Emperor's New Mind*. Oxford: Oxford University Press.

1163 Penrose, R. (2000). *The Large, the Small and the Human Mind*. Cambridge: Cambridge University
1164     Press.

1165 Perkel, D. H., and Bullock, T. H. (1968). Neural coding. *Neurosci. Res. Program Bull.* 6, 219–349.

1166 Phillips, C. G., Zeki, S., and Barlow, H. B. (1984). Localization of function in the cerebral cortex:
1167     past, present and future. *Brain* 107, 328–361.

1168 Pillow, J. W., Shlens, J., Paninski, L., Sher, A., Litke, A. M., Chichilnisky, E. J., and Simoncelli, E.
1169     P. (2008). Spatio-temporal correlations and visual signalling in a complete neuronal
1170     population. *Nature* 454, 995–999.

1171 Place, U. T. (1956). Is consciousness a brain process? *Br. J. Psychol.* 47, 44–50.

1172 Pockett, S. (2002). Difficulties with the electromagnetic field theory of consciousness. *J. Conscious.*
1173     *Stud.* 9, 51–56.

1174 Popper, K. R., and Eccles, J. C. (1977). *The self and its brain*. Berlin: Springer International.

1175 Quiroga, R. Q., and Panzeri, S. (2009). Extracting information from neuronal populations:
1176     information theory and decoding approaches. *Nat. Rev. Neurosci.* 10, 173–185.

1177 Reinagel, P., and Reid, R. C. (2000). Temporal Coding of Visual Information in the Thalamus. *J.*
1178     *Neurosci.* 20, 5392–5400.

1179 Rieke, F., de Ruyter van Steveninck, R., Warland, D., and Bialek, W. (1997). *Spikes: exploring the*
1180     *neural code*. Cambridge, MA: MIT Press.

1181 Rodriguez, E., George, N., Lachaux, J.-P., Martinerie, J., Renault, B., and Varela, F. J. (1999).
1182     Perception's shadow: long-distance synchronization of human brain activity. *Nature* 397,
1183     430–433.

1184 Rosseinsky, N. M. (2014a). Multi-subject settings restrict and define physically-orthodox neural
1185     codes for consciousness. Manuscript in Submission.

1186 Rosseinsky, N. M. (2014b). Topographic codes for consciousness are constrained to a small and
1187     distinctive set by a physically-orthodox setting. Manuscript in Submission.

1188 Rumelhart, D. E., and McClelland, J. L. (1988). *Parallel Distributed Processing: Explorations in the*
1189     *Microstructure of Cognition*. Vols. 1 and 2. Cambridge, MA: Bradford.

1190 Sejnowski, T. J., and Paulsen, O. (2006). Network oscillations: emerging computational principles. *J.*
1191     *Neurosci.* 26, 1673–1676.

1192 Shulman, R. G., Hyder, F., and Rothman, D. L. (2009). Baseline brain energy supports the state of

1193            consciousness. *Proc. Natl. Acad. Sci.* 106, 11096–11101.

1194    Smart, J. J. C. (1959). Sensations and brain processes. *Philos. Rev.* 68, 141–156.

1195    Solomon, S. G., and Lennie, P. (2007). The machinery of colour vision. *Nat. Rev. Neurosci.* 8, 276–
1196            286.

1197    Stapp, H. P. (1993). *Mind, matter, and quantum mechanics*. Berlin: Springer.

1198    Tegmark, M. (2000). Importance of quantum decoherence in brain processes. *Phys. Rev. E* 61, 4194.

1199    Thorpe, S. J., and Imbert, M. (1989). "Biological constraints on connectionist modelling," in
1200            *Connectionism in perspective*, 63–92.

1201    Tononi, G. (2004). An information integration theory of consciousness. *BMC Neurosci.* 5, 42.

1202    Tononi, G., and Edelman, G. M. (1998). Consciousness and complexity. *Science* 282, 1846–1851.

1203    Tononi, G., and Koch, C. (2008). The neural correlates of consciousness. *Ann. N. Y. Acad. Sci.* 1124,
1204            239–261.

1205    Treisman, A. (1999). Solutions to the binding problem: progress through controversy and
1206            convergence. *Neuron* 24, 105–110.

1207    Tsodyks, M. V., Skaggs, W. E., Sejnowski, T. J., and McNaughton, B. L. (1996). Population
1208            dynamics and theta rhythm phase precession of hippocampal place cell firing: a spiking
1209            neuron model. *Hippocampus* 6, 271–280.

1210    VanRullen, R., and Koch, C. (2003). Is perception discrete or continuous? *Trends Cogn. Sci.* 7, 207–
1211            213.

1212    Van Vreeswijk, C. (2004). "What is the neural code," in *23 Problems in Systems Neuroscience*, eds.
1213            J. L. van Hemmen and T. J. Sejnowski (New York: Oxford University Press), 143–159.

1214    Wagner, M. (2006). *The Geometries of Visual Space*. New York: Psychology Press.

1215    Wilkes, K. V. (1988). "---, yishi, , duh, um, and consciousness," in *Consciousness and Contemporary*
1216            *Science*, eds. A. Marcel and E. Bisiach (Oxford: Clarendon Press), 16–41.

1217    Wilkes, K. V. (1984). Is consciousness important? *Br. J. Philos. Sci.* 35, 223–243.

1218    Zaidi, Q., Victor, J., McDermott, J., Geffen, M., Bensmaia, S., and Cleland, T. A. (2013). Perceptual
1219            Spaces: Mathematical Structures to Neural Mechanisms. *J. Neurosci.* 33, 17597–17602.
1220            doi:10.1523/jneurosci.3343-13.2013.

1221
1222

1223 **10.** **Figures**

1224 **Figure 1. Schematic illustration of the primary result.** **(A)** Distributed vector encoding *e.g.* of
1225 color *as a component-of-consciousness* leads to non-locality under physical orthodoxy, because a
1226 property (conscious experience of color) depends on dynamical activity at more than one location. In
1227 a physically-orthodox theory-of-consciousness there can be no intermediate, dedicated, physical
1228 structure connecting encoding locations to conscious-experiential locations. (Double-headed black
1229 arrow schematically indicates that the red component-of-consciousness is a property of or
1230 synonymous with encoding dynamics at multiple locations.) **(B)** Distributed vector encoding *e.g.* of
1231 color *as part of conventional computational neuroscience* does not lead to non-locality, because later
1232 computational sites are connected to each vector component via intermediate physical structure (*e.g.*
1233 axons, synapses, and dendrites). The frontal spike shown, encoding *e.g.* "existence of red in the
1234 visual field", is generated by a feedforward computation that is local at each point. (Single-headed
1235 black arrow schematically depicts feedforward axon-synapse-dendrite connectivity between encoding
1236 dynamics at multiple locations and "behavioral computation" at a single location.)



1237
1238

Nicholas M. Rosseinsky

1239   **Figure 2. Schematic illustration of "parallel" and "distributed" nomenclature.** Three distinct
1240   locations in the visual field are schematically identified by colored crosses. Two topographic maps in
1241   cortical areas "1" and "2" encode different features of the visual field (oriented edges and motion
1242   respectively, say). Colored cortical locations denote neurons in each cortical area dedicated to
1243   encoding the correspondingly-colored location in the visual environment. Specifically, *one* neuron
1244   per environmental location is used in area 1, and *three* neurons per location in area 2. In this
1245   encoding scheme, feature encoding is "parallel" (in two areas), and encodings of distinct locations
1246   are computed "in parallel" (in dedicated neurons or groups of neurons, arranged topographically).
1247   Area 1 uses "non-distributed" encoding of feature instances, whereas area 2 employs a "distributed"
1248   code.



1249
1250

1251 **Figure 3. Schematic depiction of locality-of-causality conceptions in multiple-space settings.** In
1252 this Figure, the light grey space (solid boundary) is identified with the orthodox space containing the
1253 physical universe, and the dark grey space (dashed boundary) with the space containing the conscious
1254 experience associated with a single brain (not to scale). **(A)** One space embedded in another: two
1255 views of the same embedding are depicted, showing that points inside the dashed boundary belong to
1256 both the dark grey (upper image) and light grey (lower image) spaces. **(B),(C)** In the embedded
1257 construction of panel A, a neural dynamic at Y can either be non-local (panel B) or local (panel C) to
1258 a component of conscious experience Z. **(D)** Two topologically-unconnected spaces: conscious-
1259 experiential space is not spatially "inside" the physical universe (although physical duality can be
1260 avoided by proposing that conscious-experiential *space* is a property of brain dynamics). **(E)** In the
1261 unconnected construction of panel D, a neural dynamic at Y cannot be (conventionally) local to a
1262 component-of-consciousness at Z, but other neural dynamics U and V can be respectively non-local
1263 and local to Y. Thus, two distinct locality conceptions exist, *i.e.* locality of Z to {U,V,Y} and locality
1264 of {U,V,Y} to each other.



1265
1266

Nicholas M. Rosseinsky

1267 **Figure 4. Schematic depiction of phenomena.** Central analyses discuss the relationship between the
1268 exteroceptive environment (schematically depicted by the line drawing of the tree), brain dynamics
1269 that encode the environment (schematically depicted by spikes), and the conscious experience of the
1270 environment (schematically depicted by the colored tree and surroundings). The dotted line
1271 delineates objective and subjective phenomena without asserting duality: a wide variety of
1272 attributions of components-of-conscious-experience to objective brain states will be accommodated
1273 by formal symbolism. The central focus will be on *physically-orthodox* attributions that avoid duality
1274 while maintaining the necessary recognition of three basic categories of phenomena: environment;
1275 brain-dynamical states encoding environment; and, conscious-experience-of-environment.



1276
1277

1278  **Figure 5. Symbols for components of the exteroceptive environment.** Symbols $s_{a,b}$ correspond to
1279  specific features of the environment that are primitives (*i.e.* irreducible informational components) of
1280  the sensory information-processing systems. For example, each oriented edge shown here
1281  corresponds to a specific symbol, $s_{edge,i}$ say. Other members (not shown) of the complete collection
1282  $\{s_{a,b}\}$ for the visual modality might include $s_{col,j}$ (color) and $s_{move,k}$ (motion). (Blue double-headed
1283  arrows in Figures 5 to 9 schematically depict relationships between natural phenomena and symbols
1284  that must be established as part of any scientific investigation.)



1285
1286

1287    **Figure 6. Coordinate symbols for locations in the exteroceptive environment.** Symbols $_X r_j$ label
1288    locations in the exteroceptive environment at which various features labeled by $s_{a,b}$ symbols (Figure
1289    5) occur.



1290
1291

1292 **Figure 7. Symbols for brain-dynamics and coordinates for brain locations.** The symbol $A$
1293 denotes an encoding-relevant measure of brain-dynamical activity. Knowledge of $A$ at the locations
1294 (red crosses) indexed by the collection of coordinates $\{_A\mathbf{r}_i\}$ is sufficient to define the complete
1295 encoding state of the brain (for behavioral computations, but not necessarily for the generation of
1296 consciousness).



1297
1298

Nicholas M. Rosseinsky

1299  **Figure 8. Symbols for components of conscious experience.** The symbol $<s_{a,b}>$ denotes the
1300  conscious experience of the exteroceptive feature labeled by $s_{a,b}$ (Figure 5). (Components of
1301  conscious experience shown here are composite combinations of edges and colors. In an $\{s_{edge}, s_{col}\}$
1302  description of the environment, conscious experience would be described in terms of collections of
1303  $<s_{edge}>$ and $<s_{col}>$ symbols. Curly parentheses $\{...\}$ always denote a set or collection, in the present
1304  series.)



1305
1306

1307  **Figure 9. Coordinate symbols for locations in conscious experience.** Symbols $\rho_k$ label relative
1308  locations within conscious-experiential space at which various components-of-consciousness labeled
1309  by $<s_{a,b}>$ symbols (Figure 8) occur.



1310
1311

1312   **Figure 10. Schematic illustration of the formal statement of a complete theory of exteroceptive**
1313   **consciousness.** Extensive definitions of symbols (Figures 4 to 9) lead to just three basic equations
1314   that can together completely state any theory-of-consciousness. (Encoding *B*-dynamics are
1315   schematically represented by black crosses. Arrows schematically depict causal relationships.) Red
1316   arrow, red equation: the existence of a stimulus instance $s_{ab}$ at an external location $_X\mathbf{r}_j$ means that *A*-
1317   dynamics at brain locations $\{_A\mathbf{r}_i\}$ must satisfy $C_{abj} = 1$. Green line, green equation: physical coupling
1318   of *A*- and *B*-dynamics means that $C_{abj} = 1$ satisfaction generates *B*-dynamics that satisfy $D_{abj} = 1$.
1319   Black line, black equation: by definition of the causal properties of *B*-states, $D_{abj} = 1$ satisfaction
1320   causes the generation of the component-of-consciousness $<s_{ab}>$ at the conscious-experiential location
1321   $\rho_j$.



$$\exists_P s_{ab}(_X\mathbf{r}_j) \Rightarrow$$
$$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$$

$$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$$
$$\Rightarrow$$
$$D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1$$

$$D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1 \Rightarrow$$
$$\exists_P <s_{ab}>(\rho_j)$$

1322
1323

1324 **Figure 11. Sub-domain and property schemes for explaining conscious-experiential space.**
1325 **(A),(B)** If a sub-domain of orthodox physical space is used as the space of conscious experience (as
1326 illustrated in Figure 3A-C), then there is only one space, which can be viewed in two ways: in panel
1327 A, neural dynamics encoding experience are displayed as the contents of space; in panel B,
1328 experience encoded by dynamics is displayed. **(C)** If both the contents of experience and the space
1329 containing that experience are properties of neural dynamics, then the space of conscious experience
1330 can be physically and topologically separate from orthodox space (as illustrated in Figure 3D-E). In
1331 this case, there are two spaces, one containing brain dynamics (as well as the brain, the environment,
1332 and the rest of the physical Universe) and another containing conscious experience. The physically-
1333 orthodox avoidance-of-duality can be achieved in this case by considering the space of conscious
1334 experience to be ontologically junior to orthodox physical space. Whereas panels A and B depict two
1335 different views of *one* space, in the property approach shown in panel C there are two distinct spaces
1336 (schematically shown as upper ellipse and lower circle).



1337
1338

1339　**Figure 12. Spatially-extended properties in a local physical theory require a spatially-extended**
1340　**measuring apparatus or other intermediate physical structure.** Green picture elements
1341　schematically depict physical structure connecting spatially-distributed physical features with
1342　spatially-extended properties. Blue double-headed arrows depict connections between physical
1343　phenomena and theoretical symbols. Red elements delineate the specific phenomenon labeled by
1344　explicitly-displayed symbols. **(A)** A collection of molecules with various momenta (illustrated by
1345　arrows of varying direction and magnitude). **(B)** Assignment of temperature $T_1$ to the collection of
1346　molecules requires a spatially-extended thermometer (green; not to scale) whose reading (red line) is
1347　the (single-location) phenomenon corresponding to $T_1$. **(C)** A collection of spatially-distributed
1348　neural spikes that together generate a component-of-consciousness (experience of oriented-edge;
1349　contents of the red square). **(D)** The label $<s_{12}>(\rho_3)$ assigned to the component-of-consciousness can
1350　only be part of a *local* physical theory if there is intermediate physically connective structure (green
1351　lines) between neural locations and the location of the component-of-consciousness. This connective
1352　structure plays the same conceptual role as the thermometer (panel B) in creating a local theory of a
1353　spatially-extended property. Connective structure [7.2.6] could be comprised of *e.g.* exotic particles
1354　in a three-dimensional setting, or orthodox particles in a minimally six-dimensional topologically
1355　connected setting (Rosseinsky, 2014b). (Other, less orthodox, possibilities for physical connection
1356　can be hypothesized.)



1357

Figure 1.JPEG

A

Conscious
experience

B

Figure 2.JPEG



Cortical Area "2"

Cortical Area "1"

Figure 3.JPEG

Figure 4.JPEG

Figure 5.JPEG

# Phenomena

# Symbols



$$\{s_{a,b}\}$$

Figure 6.JPEG
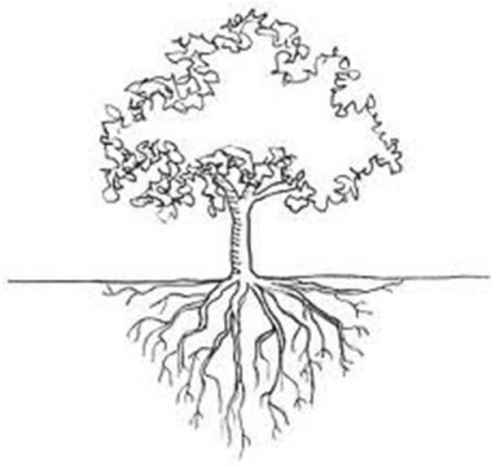
# Phenomena

# Symbols



$\{_X\mathbf{r}_j\}$

Figure 7.JPEG

# Phenomena

# Symbols

$$\{A(_A\mathbf{r}_i)\}$$

$$\{_A\mathbf{r}_i\}$$

Figure 8.JPEG

# Phenomena

# Symbols

$$\{<s_{a,b}>\}$$

Figure 9.JPEG

# Phenomena

# Symbols

$\{\rho_k\}$

Figure 10.JPEG



$$\exists_P\, s_{ab}\, (_X\mathbf{r}_j) \Rightarrow$$
$$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$$

$$C_{abj}[\{A(_A\mathbf{r}_i)\}] = 1$$
$$\Rightarrow$$
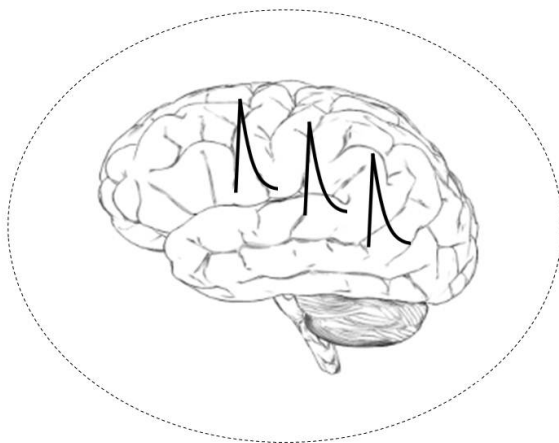$$D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1$$

$$D_{abj}[\{B(_B\mathbf{r}_i)\}] = 1 \Rightarrow$$
$$\exists_P\, <s_{ab}>(\rho_j)$$
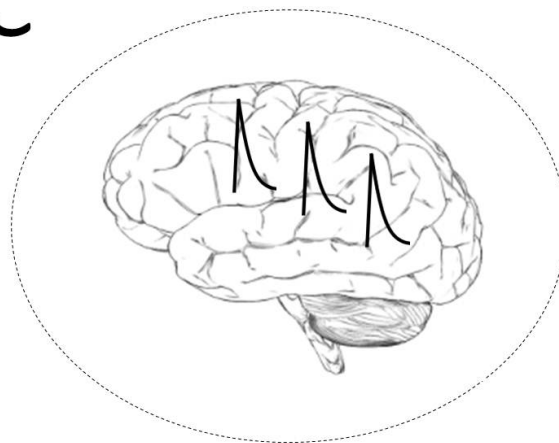
Figure 11.JPEG

Figure 12.JPEG



A

B

$T_1$

C

D

$\langle s_{1,2}\rangle(\rho_3)$