

Version 2

Machine Hyperconsciousness

Rodrick Wallace, Ph.D.
The New York State Psychiatric Institute*

July 12, 2006

Abstract

Individual animal consciousness appears limited to a single giant component of interacting cognitive modules, instantiating a shifting, highly tunable, Global Workspace. Human institutions, by contrast, can support several, often many, such giant components simultaneously, although they generally function far more slowly than the minds of the individuals who compose them. Machines having multiple global workspaces – hyperconscious machines – should, however, be able to operate at the few hundred milliseconds characteristic of individual consciousness. Such multitasking – machine or institutional – while clearly limiting the phenomenon of inattentional blindness, does not eliminate it, and introduces characteristic failure modes involving the distortion of information sent between global workspaces. This suggests that machines explicitly designed along these principles, while highly efficient at certain sets of tasks, remain subject to canonical and idiosyncratic failure patterns analogous to, but more complicated than, those explored in Wallace (2006a). By contrast, institutions, facing similar challenges, are usually deeply embedded in a highly stabilizing cultural matrix of law, custom, and tradition which has evolved over many centuries. Parallel development of analogous engineering strategies - directed toward ensuring an ‘ethical’ device – would seem requisite to the successful application of any form of hyperconscious machine technology.

Key words bandpass, cognition, consciousness, directed homotopy, global workspace, groupoid, institution, information theory, multitasking, random network, rate distortion, topology.

INTRODUCTION

For nearly a half-million years, humans and other hominids in small, well trained, disciplined groups, have been the most efficient and fearsome predators on Earth. More recently humans, in large-scale organization, have recast the surface features and ecological dynamics of the entire planet. Human organizations, at all scales, are cognitive, taking the perspective of Atlan and Cohen (1998), in that they perceive patterns of threat or opportunity, compare those patterns with

some internal, learned or inherited, picture of the world, and then choose one or a small number of responses from a vastly larger repertory of what is possible to them. Human institutions are now the subject of intense study from the perspectives of distributed cognition (e.g. Patel, 1998; Cohen et al., 2006; Laxmisan et al., 2006, and references therein). Wallace (2006b) has, in fact, attempted to apply a generalization of a currently popular theory of individual consciousness to their analysis.

Indeed, consciousness study has again become academically respectable, after nearly a century of ideologically enforced unpopularity, and Baars’ Global Workspace Theory (GWT), (Baars, 1988, 2005; Baars and Franklin, 2003) is the present front runner in the Darwinian competition between theoretical approaches (Dehaene and Naccache, 2001). Wallace and colleagues (Wallace, 2005a, b, 2006a; Glazebrook, 2006) have developed a fairly detailed mathematical model of GWT, using a Dretske-like information theory formalism (Dretske, 1981, 1988, 1993, 1994), extended by tools from statistical physics, the Large Deviations Program of applied probability, and the topological theory of highly parallel computation.

Institutional cognition involving simultaneous, multiple global workspaces is, however, a far more complex and varied phenomenon, significantly less constrained by biological evolution, and far more efficient in many important respects (Wallace, 2006b).

As the cultural anthropologists will attest, the structures, functions, and innate character of institutional cognition are greatly variable and highly adaptable across social and physical geography, and across history. Individual human consciousness, by contrast, remains constrained by the primary biological necessity of single-tasking, leading to the striking phenomenon of inattentional blindness (IAB) when the Rate Distortion Manifold of consciousness become necessarily focused on one primary process to the virtual exclusion of others which might be expected to intrude (e.g. Mack, 1998; Dehaene and Changeux, 2005; Matsuda and Nisbett, 2006; Simon and Chabris, 1999; Simons, 2000; Wayand et al., 2005).

Generalizing a second order mathematical treatment of Baars’ Global Workspace model of individual consciousness to organizational structures, Wallace (2006b) has suggested the contrasting possibility of collective multitasking, although that is mathematically a far more complicated process to an-

*Address correspondence to R. Wallace, PISCS Inc., 549 W. 123 St., Suite 16F, New York, NY, 10027. Telephone (212) 865-4766, email rd-wall@ix.netcom.com. Affiliations are for identification only.

alyze and describe. In that work he uncovered an institutional analog to individual inattentional blindness, and additional failure modes specific to the complication of multiple workspaces.

Here we will recapitulate that work from the perspective of machine design, specifically invoking machines having as many global workspaces as a large, highly capable human institution, but operating with characteristic times similar to individual consciousness – a few hundred milliseconds.

By contrast, small, disciplined groups involved in hunting, combat, firefighting, sports, or emergency medicine, may function in realms of a few seconds to an hour. The appropriate time constant probably follows something like a logarithmic scaling, i.e. $t \propto \log[N]$, where t is the time and N the number of individuals or, more likely, workspaces, which must intercoordinate across the organization.

Such fast, multiple-workspace, machines could be expected to do far more than just play variants of chess well, although their canonical and idiosyncratic failure modes would not be stabilized by the tens of thousands of years of cultural and ‘market’ selection pressures which have come to structure the various cognitive human institutions.

We start with a review of recent work on individual consciousness, as a kind of second order iteration of simple cognition, and then begin to examine the nontrivial extensions needed to describe machines having multiple workspaces.

INTRODUCTION TO THE FORMAL THEORY

1. The Global Workspace model of individual consciousness

The central ideas of Baars’ Global Workspace Theory of individual consciousness are as follows (Baars and Franklin, 2003):

- (1) The brain can be viewed as a collection of distributed specialized networks (processors).
- (2) Consciousness is associated with a global workspace in the brain – a fleeting memory capacity whose focal contents are widely distributed (broadcast) to many unconscious specialized networks.
- (3) Conversely, a global workspace can also serve to integrate many competing and cooperating input networks.
- (4) Some unconscious networks, called contexts, shape conscious contents, for example unconscious parietal maps modulate visual feature cells that underlie the perception of color in the ventral stream.
- (5) Such contexts work together jointly to constrain conscious events.
- (6) Motives and emotions can be viewed as goal contexts.
- (7) Executive functions work as hierarchies of goal contexts.

Although this basic approach has been the focus of work by many researchers for two decades, consciousness studies has only recently, in the context of a deluge of empirical results from brain imaging experiments, begun digesting the perspective and preparing to move on.

Currently popular agent-based and artificial neural network (ANN) treatments of cognition, consciousness and other higher order mental functions, to take Krebs’ (2005) view, are

little more than sufficiency arguments, in the same sense that a Fourier series expansion can be empirically fitted to nearly any function over a fixed interval without providing real understanding of the underlying structure. Necessary conditions, as Dretske argues (Dretske, 1981, 1988, 1993, 1994), give considerably more insight.

Wallace (2005a, b) addresses Baars’ theme from Dretske’s viewpoint, examining the necessary conditions which the asymptotic limit theorems of information theory impose on the Global Workspace. A central outcome of that work is the incorporation, in a natural manner, of constraints on individual consciousness, i.e. what Baars calls contexts. Using information theory methods, extended by an obvious homology between information source uncertainty and free energy density, it is possible to formally account for the effects on individual consciousness of parallel physiological modules like the immune system, embedding structures like the local social network, and, most importantly, the all-encompassing cultural heritage which so uniquely marks human biology (e.g. Richerson and Boyd, 2004). This embedding evades the mereological fallacy which fatally bedevils brain-only theories of human consciousness (Bennett and Hacker, 2003).

Transfer of phase change approaches from statistical physics to information theory via the same homology generates the punctuated nature of accession to consciousness in a similarly natural manner. The necessary renormalization calculation focuses on a phase transition driven by variation in the average strength of nondisjunctive weak ties (Granovetter, 1973) linking unconscious cognitive submodules. A second-order universality class tuning allows for adaptation of conscious attention via rate distortion manifolds which generalize the idea of a retina. The Baars model emerges as an almost exact parallel to hierarchical regression, based, however, on the Shannon-McMillan rather than the Central Limit Theorem.

Wallace (2005b) recently proposed a somewhat different approach, using classic results from random and semirandom network theory (Erdos and Renyi, 1960; Albert and Barabasi, 2002; Newman, 2003) applied to a modular network of cognitive processors. The unconscious modular network structure of the brain is, of course, not random. However, in the spirit of the wag who said “all mathematical models are wrong, but some are useful”, the method serves as the foundation of a different, but roughly parallel, treatment of the Global Workspace to that given in Wallace (2005a), and hence as another basis for a benchmark model against which empirical data can be compared.

The first step is to argue for the existence of a network of loosely linked cognitive unconscious modules, and to characterize each of them by the richness of the canonical language – information source – associated with it. This is in some contrast to attempts to explicitly model neural structures themselves using network theory, e.g. the neuropercolation approach of Kozma et al. (2004, 2005), which nonetheless uses many similar mathematical techniques. Here, rather, we look at the necessary conditions imposed by the asymptotic limits of information theory on any realization of a cognitive pro-

cess, be it biological wetware, silicon dryware, or some direct or systems-level hybrid. All cognitive processes, in this formulation, are to be associated with a canonical dual information source which will be constrained by the Rate Distortion Theorem, or, in the zero-error limit, the Shannon-McMillan Theorem, both of which are described further in the Mathematical Appendix. It is interactions between nodes in this abstractly defined network which will be of interest here, rather than whatever mechanisms, social or biological system, or mixture of them, actually constitute the underlying cognitive modules.

The second step is to examine the conditions under which a giant component (GC) suddenly emerges as a kind of phase transition in a network of such linked cognitive modules, to determine how large that component is, and to define the relation between the size of the component and the richness of the cognitive language associated with it. This is the candidate for Baars' shifting Global Workspace of consciousness.

While Wallace (2005a) examines the effect of changing the average strength of nondisjunctive weak ties acting across linked unconscious modules, Wallace (2005b) focuses on changing the average *number* of such ties having a fixed strength, a complementary perspective whose extension via a kind of 'renormalization' leads to a far more general approach.

The third step, following Wallace (2005b), is to tune the threshold at which the giant component comes into being, and to tune vigilance, the threshold for accession to consciousness.

Wallace's (2005b) information theory modular network treatment can be enriched by introducing a groupoid formalism which is roughly similar to recent analyses of linked dynamic networks described by differential equation models (e.g. Golubitsky and Stewart, 2006; Stewart et al., 2003, Stewart, 2004; Weinstein, 1996; Connes, 1994; Bak et al., 2006). Internal and external linkages between information sources break the underlying groupoid symmetry, and introduce more structure, the global workspace and the effect of contexts, respectively. The analysis provides a foundation for further mathematical exploration of linked cognitive processes.

The generalization of interest here is to examine the conditions under which cognitive modules may multitask, engaging in more than one giant component at the same time, i.e. synchronously. This is something which individual consciousness does not permit under normal circumstances. The obvious tradeoff, of course, is the very rapid flow of individual consciousness, a matter of a few hundred milliseconds, as opposed to the much slower, if considerably more comprehensive, operations of institutional cognition.

The conjecture, of course, is that machines can be built which would carry out complex processes of multiple global workspace cognition at rates approaching those of individual conscious animals.

2. Cognition as language

Cognition is not consciousness. Most mental, and many physiological, functions, while cognitive in a formal sense, hardly ever become entrained into the Global Workspace of individual consciousness: one seldom is able to consciously regulate immune function, blood pressure, or the details of

binocular tracking and bipedal motion, except to decide 'what shall I look at', 'where shall I walk'. Nonetheless, many cognitive processes, conscious or unconscious, appear intimately related to language, broadly speaking. The construction is fairly straightforward (Wallace, 2000, 2005a, b).

Atlan and Cohen (1998) and Cohen (2000) argue, in the context of immune cognition, that the essence of cognitive function involves comparison of a perceived signal with an internal, learned picture of the world, and then, upon that comparison, choice of one response from a much larger repertoire of possible responses.

Cognitive pattern recognition-and-response proceeds by an algorithmic combination of an incoming external sensory signal with an internal ongoing activity – incorporating the learned picture of the world – and triggering an appropriate action based on a decision that the pattern of sensory activity requires a response.

More formally, a pattern of sensory input is mixed in an unspecified but systematic algorithmic manner with a pattern of internal ongoing activity to create a path of combined signals $x = (a_0, a_1, \dots, a_n, \dots)$. Each a_k thus represents some functional composition of internal and external signals. Wallace (2005a) provides two neural network examples.

This path is fed into a highly nonlinear, but otherwise similarly unspecified, decision oscillator, h , which generates an output $h(x)$ that is an element of one of two disjoint sets B_0 and B_1 of possible system responses. Let

$$B_0 \equiv b_0, \dots, b_k,$$

$$B_1 \equiv b_{k+1}, \dots, b_m.$$

Assume a graded response, supposing that if

$$h(x) \in B_0,$$

the pattern is not recognized, and if

$$h(x) \in B_1,$$

the pattern is recognized, and some action $b_j, k+1 \leq j \leq m$ takes place.

The principal objects of formal interest are paths x which trigger pattern recognition-and-response. That is, given a fixed initial state a_0 , we examine all possible subsequent paths x beginning with a_0 and leading to the event $h(x) \in B_1$. Thus $h(a_0, \dots, a_j) \in B_0$ for all $0 < j < m$, but $h(a_0, \dots, a_m) \in B_1$.

For each positive integer n , let $N(n)$ be the number of high probability grammatical and syntactical paths of length n which begin with some particular a_0 and lead to the condition $h(x) \in B_1$. Call such paths 'meaningful', assuming, not unreasonably, that $N(n)$ will be considerably less than the number of all possible paths of length n leading from a_0 to the condition $h(x) \in B_1$.

While combining algorithm, the form of the nonlinear oscillator, and the details of grammar and syntax, are all unspecified in this model, the critical assumption which permits

inference on necessary conditions constrained by the asymptotic limit theorems of information theory is that the finite limit

$$(1) \quad H \equiv \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n}$$

both exists and is independent of the path x .

We call such a pattern recognition-and-response cognitive process *ergodic*. Not all cognitive processes are likely to be ergodic, implying that H , if it indeed exists at all, is path dependent, although extension to nearly ergodic processes, in a certain sense, seems possible (Wallace, 2005a).

Invoking the spirit of the Shannon-McMillan Theorem, it is possible to define an adiabatically, piecewise stationary, ergodic information source \mathbf{X} associated with stochastic variates X_j having joint and conditional probabilities $P(a_0, \dots, a_n)$ and $P(a_n|a_0, \dots, a_{n-1})$ such that appropriate joint and conditional Shannon uncertainties satisfy the classic relations

$$H[\mathbf{X}] = \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n} =$$

$$\lim_{n \rightarrow \infty} H(X_n|X_0, \dots, X_{n-1}) =$$

$$\lim_{n \rightarrow \infty} \frac{H(X_0, \dots, X_n)}{n}.$$

This information source is defined as *dual* to the underlying ergodic cognitive process (Wallace, 2005a).

Recall that the Shannon uncertainties $H(\dots)$ are cross-sectional law-of-large-numbers sums of the form $-\sum_k P_k \log[P_k]$, where the P_k constitute a probability distribution. See Khinchin (1957), Ash (1990), or Cover and Thomas (1991) for the standard details.

3. The cognitive modular network symmetry groupoid

A formal equivalence class algebra can be constructed by choosing different origin points a_0 and defining equivalence by the existence of a high probability meaningful path connecting two points. Disjoint partition by equivalence class, analogous to orbit equivalence classes for dynamical systems, defines the vertices of the proposed network of cognitive dual languages. Each vertex then represents a different information source dual to a cognitive process. This is not a representation of a neural network as such, or of some circuit in silicon. It is, rather, an abstract set of ‘languages’ dual to the cognitive processes instantiated by biological structures, social process, machines, or their hybrids.

This structure is a groupoid, in the sense of Weinstein (1996). States a_j, a_k in a set A are related by the groupoid

morphism if and only if there exists a high probability grammatical path connecting them, and tuning across the various possible ways in which that can happen – the different cognitive languages – parametrizes the set of equivalence relations and creates the groupoid. This assertion requires some development.

Note that not all possible pairs of states (a_j, a_k) can be connected by such a morphism, i.e. by a high probability, grammatical and syntactical cognitive path, but those that can define the groupoid element, a morphism $g = (a_j, a_k)$ having the natural inverse $g^{-1} = (a_k, a_j)$. Given such a pairing, connection by a meaningful path, it is possible to define ‘natural’ end-point maps $\alpha(g) = a_j, \beta(g) = a_k$ from the set of morphisms G into A , and a formally associative product in the groupoid $g_1 g_2$ provided $\alpha(g_1 g_2) = \alpha(g_1), \beta(g_1 g_2) = \beta(g_2)$, and $\beta(g_1) = \alpha(g_2)$. Then the product is defined, and associative, i.e. $(g_1 g_2) g_3 = g_1 (g_2 g_3)$.

In addition there are natural left and right identity elements λ_g, ρ_g such that $\lambda_g g = g = g \rho_g$ whose characterization is left as an exercise (Weinstein, 1996).

An orbit of the groupoid G over A is an equivalence class for the relation $a_j \sim G a_k$ if and only if there is a groupoid element g with $\alpha(g) = a_j$ and $\beta(g) = a_k$.

The isotopy group of $a \in X$ consists of those g in G with $\alpha(g) = a = \beta(g)$.

In essence a groupoid is a category in which all morphisms have an inverse, here defined in terms of connection by a meaningful path of an information source dual to a cognitive process.

If G is any groupoid over A , the map $(\alpha, \beta) : G \rightarrow A \times A$ is a morphism from G to the pair groupoid of A . The image of (α, β) is the orbit equivalence relation $\sim G$, and the functional kernel is the union of the isotropy groups. If $f : X \rightarrow Y$ is a function, then the kernel of f , $\ker(f) = [(x_1, x_2) \in X \times X : f(x_1) = f(x_2)]$ defines an equivalence relation.

As Weinstein (1996) points out, the morphism (α, β) suggests another way of looking at groupoids. A groupoid over A identifies not only which elements of A are equivalent to one another (isomorphic), but *it also parametrizes the different ways (isomorphisms) in which two elements can be equivalent*, i.e. all possible information sources dual to some cognitive process. Given the information theoretic characterization of cognition presented above, this produces a full modular cognitive network in a highly natural manner.

The groupoid approach has become quite popular in the study of networks of coupled dynamical systems which can be defined by differential equation models, (e.g. Golubitsky and Stewart, 2006; Stewart et al. (2003), Stewart (2004)). Here we have outlined how to extend the technique to networks of interacting information sources which, in a dual sense, characterize cognitive processes, and cannot at all be described by the usual differential equation models. These latter, it seems, are much the spiritual offspring of 18th Century mechanical clock models. Cognitive and conscious processes in humans involve neither computers nor clocks, but remain constrained by the limit theorems of information theory, and these permit scientific inference on necessary conditions.

4. Internal forces breaking the symmetry groupoid

The symmetry groupoid, as we have constructed it for cognitive modules, in a kind of information space, is parametrized across that space by the possible ways in which states a_j, a_k can be equivalent, i.e. connected by a meaningful path of an information source dual to a cognitive process. These are different, and in this approximation, non-interacting cognitive processes. But symmetry groupoids, like symmetry groups, are made to be broken: by internal cross-talk akin to spin-orbit interactions within a symmetric atom, and by cross-talk with slower, external, information sources, akin to putting a symmetric atom in a powerful magnetic or electric field.

As to the first process, suppose that linkages can fleetingly occur between the ordinarily disjoint cognitive modules defined by the network groupoid. In the spirit of Wallace (2005a), this is represented by establishment of a non-zero mutual information measure between them: a cross-talk which breaks the strict groupoid symmetry developed above.

Wallace (2005a) describes this structure in terms of fixed magnitude disjunctive strong ties which give the equivalence class partitioning of modules, and nondisjunctive weak ties which link modules across the partition, and parametrizes the overall structure by the average strength of the weak ties, to use Granovetter's (1973) term. By contrast the approach of Wallace (2005b), which we outline here, is to simply look at the average number of fixed-strength nondisjunctive links in a random topology. These are obviously the two analytically tractable limits of a much more complicated regime.

Since we know nothing about how the cross-talk connections can occur, we will – at first – assume they are random and construct a random graph in the classic Erdos/Renyi manner. Suppose there are M disjoint cognitive modules – M elements of the equivalence class algebra of languages dual to some cognitive process – which we now take to be the vertices of a possible graph.

For M very large, following Savante et al. (1993), when edges (defined by establishment of a fixed-strength mutual information measure between the graph vertices) are added at random to M initially disconnected vertices, a remarkable transition occurs when the number of edges becomes approximately $M/2$. Erdos and Renyi (1960) studied random graphs with M vertices and $(M/2)(1 + \mu)$ edges as $M \rightarrow \infty$, and discovered that such graphs almost surely have the following properties (Molloy and Reed, 1995, 1998; Grimmett and Stacey, 1998; Luczak, 1990; Aiello et al., 200; Albert and Barabasi, 2002):

[1] If $\mu < 0$, only small trees and unicyclic components are present, where a unicyclic component is a tree with one additional edge; moreover, the size of the largest tree component is $(\mu - \ln(1 + \mu))^{-1} + \mathcal{O}(\log \log n)$.

[2] If $\mu = 0$, however, the largest component has size of order $M^{2/3}$.

[3] If $\mu > 0$, there is a unique giant component (GC) whose size is of order M ; in fact, the size of this component is asymptotically αM , where $\mu = -\alpha^{-1}[\ln(1 - \alpha) - 1]$, which has an explicit solution for α in terms of the Lambert W-function. Thus, for example, a random graph with approxi-

mately $M \ln(2)$ edges will have a giant component containing $\approx M/2$ vertices.

Such a phase transition initiates a new, collective, cognitive phenomenon. At the level of the individual mind, unconscious cognitive modules link up to become the Global Workspace of consciousness, emergently defined by a set of cross-talk mutual information measures between interacting unconscious cognitive submodules. The source uncertainty, H , of the language dual to the collective cognitive process, which characterizes the richness of the cognitive language of the workspace, will grow as some monotonic function of the size of the GC, as more and more unconscious processes are incorporated into it. Wallace (2005b) provides details.

Others have taken similar network phase transition approaches to assemblies of neurons, e.g. neuropercolation (Kozma et al., 2004, 2005), but their work has not focused explicitly on modular networks of cognitive processes, which may or may not be instantiated by neurons. Restricting analysis to such modular networks finesses much of the underlying conceptual difficulty, and permits use of the asymptotic limit theorems of information theory and the import of techniques from statistical physics, a matter we will discuss later.

5. External forces breaking the symmetry groupoid

Just as a higher order information source, associated with the GC of a random or semirandom graph, can be constructed out of the interlinking of unconscious cognitive modules by mutual information, so too external information sources, for example in humans the cognitive immune and other physiological systems, and embedding sociocultural structures, can be represented as slower-acting information sources whose influence on the GC can be felt in a collective mutual information measure. For machines or institutions these would be the onion-like 'structured environment', to be viewed as among Baars' contexts (Baars, 1988, 2005; Baars and Franklin, 2003). The collective mutual information measure will, through the Joint Asymptotic Equipartition Theorem which generalizes the Shannon-McMillan Theorem, be the splitting criterion for high and low probability joint paths across the entire system.

The tool for this is network information theory (Cover and Thomas, 1991, p. 388). Given three interacting information sources, Y_1, Y_2, Z , the splitting criterion, taking Z as the 'external context', is given by

$$I(Y_1, Y_2|Z) = H(Z) + H(Y_1|Z) + H(Y_2|Z) - H(Y_1, Y_2, Z), \quad (2)$$

where $H(..|..)$ and $H(.,.,..)$ represent conditional and joint uncertainties (Khinchin, 1957; Ash, 1990; Cover and Thomas, 1991).

This generalizes to

$$I(Y_1, \dots, Y_n | Z) = H(Z) + \sum_{j=1}^n H(Y_j | Z) - H(Y_1, \dots, Y_n, Z).$$

(3)

If we assume the Global Workspace/Giant Component to involve a very rapidly shifting, and indeed highly tunable, dual information source X , embedding contextual cognitive modules like the immune system will have a set of significantly slower-responding sources $Y_j, j = 1..m$, and external social, cultural and other environmental processes will be characterized by even more slowly-acting sources $Z_k, k = 1..n$. Mathematical induction on equation (3) gives a complicated expression for a mutual information splitting criterion which we write as

$$I(X | Y_1, \dots, Y_m | Z_1, \dots, Z_n).$$

(4)

This encompasses a fully interpenetrating biopsychosociocultural structure for individual consciousness, one in which Baars' contexts act as important, but flexible, boundary conditions, defining the underlying topology available to the far more rapidly shifting global workspace (Wallace, 2005a, b).

This result does not commit the mereological fallacy which Bennett and Hacker (2003) impute to excessively neurocentric perspectives on consciousness in humans, that is, the mistake of imputing to a part of a system the characteristics which require functional entirety. The underlying concept of this fallacy should extend to machines interacting with their environments, and its baleful influence probably accounts for a significant part of the failure of Artificial Intelligence to deliver. See Wallace (2006) for further discussion.

6. Punctuation phenomena

As a number of researchers have noted, in one way or another, – see Wallace, (2005a) for discussion – equation (1),

$$H \equiv \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n},$$

is homologous to the thermodynamic limit in the definition of the free energy density of a physical system. This has the form

$$F(K) = \lim_{V \rightarrow \infty} \frac{\log[Z(K)]}{V},$$

(5)

where F is the free energy density, K the inverse temperature, V the system volume, and $Z(K)$ is the partition function defined by the system Hamiltonian.

Wallace (2005a) shows at some length how this homology permits the natural transfer of renormalization methods from statistical mechanics to information theory. In the spirit of the Large Deviations Program of applied probability theory, this produces phase transitions and analogs to evolutionary punctuation in systems characterized by piecewise, adiabatically stationary, ergodic information sources. These biological phase changes appear to be ubiquitous in natural systems and can be expected to dominate machine behaviors as well, particularly those which seek to emulate biological paradigms. Wallace (2002) uses these arguments to explore the differences and similarities between evolutionary punctuation in genetic and learning plateaus in neural systems.

7. Multiple Workspaces

The random network development above is predicated on there being a variable average number of fixed-strength linkages between components. Clearly, the mutual information measure of cross-talk is not inherently fixed, but can continuously vary in magnitude. This we address by a parametrized renormalization. In essence the modular network structure linked by mutual information interactions has a topology depending on the degree of interaction of interest. Suppose we define an interaction parameter ω , a real positive number, and look at geometric structures defined in terms of linkages which are zero if mutual information is less than, and 'renormalized' to unity if greater than, ω . Any given ω will define a regime of giant components of network elements linked by mutual information greater than or equal to it.

The fundamental conceptual trick at this point is to invert the argument: A given topology for the giant component will, in turn, define some critical value, ω_C , so that network elements interacting by mutual information less than that value will be unable to participate, i.e. will be locked out and not be consciously perceived. We hence are assuming that the ω is a tunable, syntactically-dependent, detection limit, and depends critically on the instantaneous topology of the giant component defining, for the human mind, the global workspace of consciousness. That topology is, fundamentally, the basic tunable syntactic filter across the underlying modular symmetry groupoid, and variation in ω is only one aspect of a much more general topological shift. More detailed analysis is given below in terms of a topological rate distortion manifold.

There is considerable empirical evidence from fMRI brain imaging experiments to show that individual human consciousness involves a single global workspace, a matter leading necessarily to the phenomenon of inattentional blindness. Cognitive submodules within institutions, – individuals, departments, formal and informal workgroups – by contrast, can do more than one thing, and indeed, are usually required to multitask. The intent of this work is to suggest the possibility of constructing machines which work on similar principles,

but much more rapidly.

Clearly multiple workspaces would lessen the probability of inattentional blindness, but, we will find, do not eliminate it, and introduce other failure modes examined in more detail later.

We must postulate a set of crosstalk information measures between cognitive submodules, each associated with its own tunable giant component having its own special topology.

Suppose the set of giant components at some ‘time’ k is characterized by a set of parameters $\Omega_k \equiv \omega_1^k, \dots, \omega_m^k$. Fixed parameter values define a particular giant component set having a particular set of topological structures. Suppose that, over a sequence of ‘times’ the set of giant components can be characterized by a (possibly coarse-grained) path $x_n = \Omega_0, \Omega_1, \dots, \Omega_{n-1}$ having significant serial correlations which, in fact, permit definition of an adiabatically, piecewise stationary, ergodic (APSE) information source in the sense of Wallace (2005a). Call that information source \mathbf{X} .

Suppose, again in the manner of Wallace (2005a), that a set of (external or internal) signals impinging on the set of giant components, is also highly structured and forms another APSE information source \mathbf{Y} which interacts not only with the system of interest globally, but specifically with the tuning parameters of the set of giant components characterized by \mathbf{X} . \mathbf{Y} is necessarily associated with a set of paths y_n .

Pair the two sets of paths into a joint path $z_n \equiv (x_n, y_n)$, and invoke some inverse coupling parameter, K , between the information sources and their paths. By the arguments of Wallace (2005a) this leads to phase transition punctuation of $I[K]$, the mutual information between \mathbf{X} and \mathbf{Y} , under either the Joint Asymptotic Equipartition Theorem, or, given a distortion measure, under the Rate Distortion Theorem.

$I[K]$ is a splitting criterion between high and low probability pairs of paths, and partakes of the homology with free energy density described in Wallace (2005a). Attentional focusing by the institution or machine then itself becomes a punctuated event in response to increasing linkage between the structure of interest and an external signal, or some particular system of internal events. This iterated argument parallels the extension of the General Linear Model into the Hierarchical Linear Model of regression theory.

Call this the Multitasking Hierarchical Cognitive Model (MHCM). For individual consciousness, there is only one giant component. For an institution, there will be a larger, and often very large, set of them. For a useful machine, the giant components must operate much more rapidly than is possible for an institution.

This requirement leads to the possibility of new failure modes related to impaired communication between Giant Components.

That is, a complication specific to high order institutional cognition or machine hyperconsciousness lies in the necessity of information transfer between giant components. The form and function of such interactions will, of course, be determined by the nature of the particular institution or machine, but, synchronous or asynchronous, contact between giant components is circumscribed by the Rate Distortion Theo-

rem. That theorem, reviewed in the Mathematical Appendix, states that, for a given maximum acceptable critical average signal distortion, there is a limiting maximum information transmission rate, such that messages sent at less than that limit are guaranteed to have average distortion less than the critical maximum. Too rapid transmission between parallel global workspaces – information overload – violates that condition, and guarantees large signal distortion. This is a likely failure mode which appears unique to multiple workspace systems which, we will argue, may otherwise have a lessened probability of inattentional blindness.

8. Cognitive quasi-thermodynamics

A fundamental homology between the information source uncertainty dual to a cognitive process and the free energy density of a physical system arises, in part, from the formal similarity between their definitions in the asymptotic limit. Information source uncertainty can be defined as in equation (1). This is quite analogous to the free energy density of a physical system, equation (5).

Feynman (1996) provides a series of physical examples, based on Bennett’s work, where this homology is, in fact, an identity, at least for very simple systems. Bennett argues, in terms of irreducibly elementary computing machines, that the information contained in a message can be viewed as the work saved by not needing to recompute what has been transmitted.

Feynman explores in some detail Bennett’s microscopic machine designed to extract useful work from a transmitted message. The essential argument is that computing, in any form, takes work, the more complicated a cognitive process, measured by its information source uncertainty, the greater its energy consumption, and our ability to provide energy to the brain is limited. Inattentional blindness, we will argue, emerges as an inevitable thermodynamic limit on processing capacity in a topologically-fixed global workspace, i.e. one which has been strongly configured about a particular task.

Understanding the time dynamics of cognitive systems away from phase transition critical points requires a phenomenology similar to the Onsager relations of nonequilibrium thermodynamics. If the dual source uncertainty of a cognitive process is parametrized by some vector of quantities $\mathbf{K} \equiv (K_1, \dots, K_m)$, then, in analogy with nonequilibrium thermodynamics, gradients in the K_j of the *disorder*, defined as

$$S \equiv H(\mathbf{K}) - \sum_{j=1}^m K_j \partial H / \partial K_j \quad (6)$$

become of central interest.

Equation (6) is similar to the definition of entropy in terms of the free energy density of a physical system, as suggested

by the homology between free energy density and information source uncertainty described above.

Pursuing the homology further, the generalized Onsager relations defining temporal dynamics become

$$dK_j/dt = \sum_i L_{j,i} \partial S / \partial K_i, \quad (7)$$

where the $L_{j,i}$ are, in first order, constants reflecting the nature of the underlying cognitive phenomena.

The L-matrix is to be viewed empirically, in the same spirit as the slope and intercept of a regression model, and may have structure far different than familiar from more simple chemical or physical processes.

The $\partial S / \partial K$ are analogous to thermodynamic forces in a chemical system, and may be subject to override by external physiological driving mechanisms (Wallace, 2005c).

Equations (6) and (7) can be derived in a simple parameter-free covariant manner which relies on the underlying topology of the information source space implicit to the development. We suppose that different physiological cognitive phenomena have, in the sense of Wallace (2000, 2005, Ch. 3), dual information sources, and are interested in the local properties of the system near a particular reference state. We impose a topology on the system, so that, near a particular ‘language’ A , dual to an underlying cognitive process, there is (in some sense) an open set U of closely similar languages \hat{A} , such that $A, \hat{A} \subset U$. Note that it may be necessary to coarse-grain the physiological responses to define these information sources. The problem is to proceed in such a way as to preserve the underlying essential topology, while eliminating ‘high frequency noise’. The formal tools for this can be found, e.g., in Chapter 8 of Burago et al. (2001).

Since the information sources dual to the cognitive processes are similar, for all pairs of languages A, \hat{A} in U , it is possible to:

- [1] Create an embedding alphabet which includes all symbols allowed to both of them.
- [2] Define an information-theoretic distortion measure in that extended, joint alphabet between any high probability (i.e. grammatical and syntactical) paths in A and \hat{A} , which we write as $d(Ax, \hat{A}x)$ (Cover and Thomas, 1991). Note that these languages do not interact, in this approximation.
- [3] Define a metric on U , for example,

$$\mathcal{M}(A, \hat{A}) = \left| \lim \frac{\int_{A, \hat{A}} d(Ax, \hat{A}x)}{\int_{A, A} d(Ax, A\hat{x})} - 1 \right|, \quad (8)$$

using an appropriate integration limit argument over the high probability paths. Note that the integration in the denominator is over different paths within A itself, while in the numerator it is between different paths in A and \hat{A} .

Consideration suggests \mathcal{M} is a formal metric, having $\mathcal{M}(A, B) \geq 0, \mathcal{M}(A, A) = 0, \mathcal{M}(A, B) = \mathcal{M}(B, A), \mathcal{M}(A, C) \leq \mathcal{M}(A, B) + \mathcal{M}(B, C)$.

Other approaches to constructing a metric on U may be possible.

Since H and \mathcal{M} are both scalars, a ‘covariant’ derivative can be defined directly as

$$dH/d\mathcal{M} = \lim_{\hat{A} \rightarrow A} \frac{H(A) - H(\hat{A})}{\mathcal{M}(A, \hat{A})}, \quad (9)$$

where $H(A)$ is the source uncertainty of language A .

Suppose the system to be set in some reference configuration A_0 .

To obtain the unperturbed dynamics of that state, we impose a Legendre transform using this derivative, defining another scalar

$$S \equiv H - \mathcal{M} dH/d\mathcal{M}. \quad (10)$$

The simplest possible Onsager relation – again an empirical equation like a regression model – in this case becomes

$$d\mathcal{M}/dt = L dS/d\mathcal{M}, \quad (11)$$

where t is the time and $dS/d\mathcal{M}$ represents an analog to the thermodynamic force in a chemical system. This is seen as acting on the reference state A_0 . For

$$dS/d\mathcal{M}|_{A_0} = 0,$$

$$d^2 S/d\mathcal{M}^2|_{A_0} > 0$$

(12)

the system is quasistable, a Black hole, if you will, and externally imposed forcing mechanisms will be needed to effect a transition to a different state.

Conversely, changing the direction of the second condition, so that

$$dS^2/d\mathcal{M}^2|_{A_0} < 0,$$

leads to a repulsive peak, a White hole, representing a possibly unattainable realm of states.

Explicit parametrization of \mathcal{M} introduces standard – and quite considerable – notational complications (e.g. Burago et al., 2001; Auslander, 1967): Imposing a metric for different cognitive dual languages parametrized by \mathbf{K} leads to Riemannian, or even Finsler, geometries (Wallace, 2005c), including the usual geodesics.

9. The simplest rate distortion manifold

The second order iteration above – analogous to expanding the General Linear Model to the Hierarchical Linear Model – which involved paths in parameter space, can itself be significantly extended. This produces a generalized tunable retina model which can be interpreted as a ‘Rate Distortion manifold’, a concept which further opens the way for import of a vast array of tools from geometry and topology.

Suppose, now, that threshold behavior for institutional reaction requires some elaborate system of nonlinear relationships defining a set of renormalization parameters $\Omega_k \equiv \omega_1^k, \dots, \omega_m^k$. The critical assumption is that there is a tunable zero order state, and that changes about that state are, in first order, relatively small, although their effects on punctuated process may not be at all small. Thus, given an initial m -dimensional vector Ω_k , the parameter vector at time $k+1$, Ω_{k+1} , can, in first order, be written as

$$\Omega_{k+1} \approx \mathbf{R}_{k+1}\Omega_k,$$

(13)

where \mathbf{R}_{t+1} is an $m \times m$ matrix, having m^2 components.

If the initial parameter vector at time $k=0$ is Ω_0 , then at time k

$$\Omega_k = \mathbf{R}_k \mathbf{R}_{k-1} \dots \mathbf{R}_1 \Omega_0.$$

(14)

The interesting correlates of individual, institutional or machine consciousness are, in this development, *now represented*

by an information-theoretic path defined by the sequence of operators \mathbf{R}_k , each member having m^2 components. The grammar and syntax of the path defined by these operators is associated with a dual information source, in the usual manner.

The effect of an information source of external signals, \mathbf{Y} , is now seen in terms of more complex joint paths in Y and R -space whose behavior is, again, governed by a mutual information splitting criterion according to the JAEPT.

The complex sequence in m^2 -dimensional R -space has, by this construction, been projected down onto a parallel path, the smaller set of m -dimensional ω -parameter vectors $\Omega_0, \dots, \Omega_k$.

If the punctuated tuning of institutional or machine attention is now characterized by a ‘higher’ dual information source – an embedding generalized language – so that the paths of the operators \mathbf{R}_k are autocorrelated, then the autocorrelated paths in Ω_k represent output of a parallel information source which is, given Rate Distortion limitations, apparently a grossly simplified, and hence highly distorted, picture of the ‘higher’ conscious process represented by the R -operators, having m as opposed to $m \times m$ components.

High levels of distortion may not necessarily be the case for such a structure, *provided it is properly tuned to the incoming signal*. If it is inappropriately tuned, however, then distortion may be extraordinary.

Let us examine a single iteration in more detail, assuming now there is a (tunable) zero reference state, \mathbf{R}_0 , for the sequence of operators \mathbf{R}_k , and that

$$\Omega_{k+1} = (\mathbf{R}_0 + \delta\mathbf{R}_{k+1})\Omega_k,$$

(15)

where $\delta\mathbf{R}_k$ is ‘small’ in some sense compared to \mathbf{R}_0 .

Note that in this analysis the operators \mathbf{R}_k are, implicitly, determined by linear regression. We thus can invoke a quasi-diagonalization in terms of \mathbf{R}_0 . Let \mathbf{Q} be the matrix of eigenvectors which Jordan-block-diagonalizes \mathbf{R}_0 . Then

$$\mathbf{Q}\Omega_{k+1} = (\mathbf{Q}\mathbf{R}_0\mathbf{Q}^{-1} + \mathbf{Q}\delta\mathbf{R}_{k+1}\mathbf{Q}^{-1})\mathbf{Q}\Omega_k.$$

(16)

If $\mathbf{Q}\Omega_k$ is an eigenvector of \mathbf{R}_0 , say Y_j with eigenvalue λ_j , it is possible to rewrite this equation as a generalized spectral expansion

$$Y_{k+1} = (\mathbf{J} + \delta\mathbf{J}_{k+1})Y_j \equiv \lambda_j Y_j + \delta Y_{k+1}$$

$$= \lambda_j Y_j + \sum_{i=1}^n a_i Y_i.$$

(17)

\mathbf{J} is a block-diagonal matrix, $\delta\mathbf{J}_{k+1} \equiv \mathbf{Q}\mathbf{R}_{k+1}\mathbf{Q}^{-1}$, and δY_{k+1} has been expanded in terms of a spectrum of the eigenvectors of \mathbf{R}_0 , with

$$|a_i| \ll |\lambda_j|, |a_{i+1}| \ll |a_i|.$$

(18)

The point is that, provided \mathbf{R}_0 has been tuned so that this condition is true, the first few terms in the spectrum of this iteration of the eigenstate will contain most of the essential information about $\delta\mathbf{R}_{k+1}$. This appears quite similar to the detection of color in the retina, where three overlapping non-orthogonal eigenmodes of response are sufficient to characterize a huge plethora of color sensation. Here, if such a tuned spectral expansion is possible, a very small number of observed eigenmodes would suffice to permit identification of a vast range of changes, so that the rate-distortion constraints become quite modest. That is, there will not be much distortion in the reduction from paths in R -space to paths in Ω -space. Inappropriate tuning, however, can produce very marked distortion, even institutional or machine inattentive blindness, in spite of multitasking.

Note that higher order Rate Distortion Manifolds are likely to give better approximations than lower ones, in the same sense that second order tangent structures give better, if more complicated, approximations in conventional differentiable manifolds (e.g. Pohl, 1962).

Indeed, Rate Distortion Manifolds can be quite formally described using standard techniques from topological manifold theory (Glazebrook, 2006). The essential point is that a rate distortion manifold is a topological structure which constrains the ‘multifactorial stream of institutional or machine consciousness’ as well as the pattern of communication between giant components, much the way a riverbank constrains the flow of the river it contains. This is a fundamental insight, which we pursue further.

10. The topology of machine cognition

The groupoid treatment of modular cognitive networks above defined equivalence classes of *states* according to whether they could be linked by grammatical/syntactical high probability ‘meaningful’ paths. Next we ask the precisely

complementary question regarding *paths*: For any two particular given states, is there some sense in which we can define equivalence classes across the set of meaningful paths linking them?

This is of particular interest to the second order hierarchical model which, in effect, describes a universality class tuning of the renormalization parameters characterizing the dancing, flowing, tunably punctuated accession to consciousness.

A closely similar question is central to recent algebraic geometry approaches to concurrent, i.e. highly parallel, computing (e.g. Pratt, 1991; Goubault and Raussen, 2002; Goubault, 2003), which we adapt.

For the moment we restrict the analysis to a giant component system characterized by two renormalization parameters, say ω_1 and ω_2 , and consider the set of meaningful paths connecting two particular points, say a and b , in the two dimensional ω -space plane of figure 1. The arguments surrounding equations (6), (7) and (12) suggests that there may be regions of fatal attraction and strong repulsion, Black holes and White holes, which can either trap or deflect the path of institutional or multitasking machine cognition.

Figures 1a and 1b show two possible configurations for a Black and a White hole, diagonal and cross-diagonal. If one requires path monotonicity – always increasing or remaining the same – then, following, e.g. Goubault (2003, figs. 6,7), there are, intuitively, two direct ways, without switchbacks, that one can get from a to b in the diagonal geometry of figure 1a, without crossing a Black or White hole, but there are three in the cross-diagonal structure of figure 1b.

Elements of each ‘way’ can be transformed into each other by continuous deformation without crossing either the Black or White hole. Figure 1a has two additional possible monotonic ways, involving over/under switchbacks, which are not drawn. Relaxing the monotonicity requirement generates a plethora of other possibilities, e.g. loopings and backwards switchbacks, whose consideration is left as an exercise. It is not clear under what circumstances such complex paths can be meaningful, a matter for further study.

These ways are the equivalence classes defining the topological structure of the two different ω -spaces, analogs to the fundamental homotopy groups in spaces which admit of loops (e.g. Lee, 2000). The closed loops needed for classical homotopy theory are impossible for this kind of system because of the ‘flow of time’ defining the output of an information source – one goes from a to b , although, for nonmonotonic paths, intermediate looping would seem possible. The theory is thus one of directed homotopy, dihomotopy, and the central question revolves around the continuous deformation of paths in ω -space into one another, without crossing Black or White holes. Goubault and Raussen (2002) provide another introduction to the formalism.

These ideas can, of course, be applied to lower level cognitive modules as well as to the second order hierarchical cognitive model of institutional or machine cognition where they are, perhaps, of more central interest.

We propose that empirical study will show how the influence of cultural heritage or developmental history defines

quite different dihomotopies of attentional focus in human organizations. That is, the topology of blind spots and their associated patterns of perceptual completion in human organizations will be culturally or developmentally modulated. It is this developmental cultural topology of multitasking organization attention which, acting in concert with the inherent limitations of the rate distortion manifold, generates the pattern of organizational inattentional blindness. Analogous developmental arguments should apply to hyperconscious machines.

Such considerations, and indeed the Black Hole development of equation (12), suggest that a multitasking organization or machine which becomes trapped in a particular pattern of behavior cannot, in general, expect to emerge from it in the absence of external forcing mechanisms.

This sort of behavior is central to ecosystem resilience theory (Gunderson, 2000; Holling, 1973), a matter which Wallace (2006b) explores in more detail.

The topology of institutional cognition provides a tool for study of resilience in human organizations or social systems, and, according to our perspective, probably for machines as well. Apparently the set of directed homotopy equivalence classes described above formally classifies quasi-equilibrium states, and thus characterizes the different possible resilience modes.

DISCUSSION AND CONCLUSIONS

The simple groupoid defined by a hyperconscious machine's basic cognitive modular structure can be broken by intrusion of (rapid) crosstalk within it, and by the imposition of (slower) crosstalk from without. The former, if strong enough, can initiate a set of topologically-determined giant component global workspaces, in a punctuated manner, while the latter deform the underlying topology of the entire system, the directed homotopy limiting what paths can actually be traversed. Broken symmetry creates richer structure in systems characterized by groupoids, just as it does for those characterized by groups.

Multitasking machine attention, in this picture, acts through a Rate Distortion manifold, a kind of retina-like filter for grammatical and syntactical meaningful paths. Signals outside the topologically constrained tunable syntax/grammar bandpass of this manifold are subject to lessened probability of punctuated conscious detection: inattentional blindness. Path-dependent machine developmental history will, according to this model, profoundly affect the phenomenon by imposing additional topological constraints defining the 'surface' along which this second order behavior can (and cannot) glide.

Glazebrook (2006) has suggested that, lurking in the background of this basic construction, is what Bak et al. (2006) call a groupoid atlas, i.e. an extension of topological manifold theory to groupoid mappings. Also lurking is identification and exploration of the natural groupoid convolution algebra which so often marks these structures (e.g. Weinstein, 1996; Connes, 1994).

Consideration suggests, in fact, that a path may be meaningful according to the groupoid parametrization of all possible

dual information sources, and that tuning is done across that parametrization via a rate distortion manifold.

Implicit, however, are the constraints imposed by machine history, in a large sense, which may further limit the properties of \mathbf{R}_0 , i.e. hold it to a developmentally determined topology.

Here we have attempted to reexpress this trade-off in terms of a syntactical/grammatical version of conventional signal theory, i.e. as a 'tuned meaningful path' form of the classic balance between sensitivity and selectivity, as particularly constrained by the directed homotopy imposed by a machine experience that is itself the outcome a historical process involving interaction with an external environment.

Overall, this analysis is analogous to, but more complicated than, Wallace's information dynamics instantiation of Baars' Global Workspace theory (Wallace, 2005a, b). Intuitively, one suspects that the higher the dimension of the second order attentional Rate Distortion Manifold, that is, the greater the multitasking, the broader the effective bandwidth of attentional focus, and the less likely is inattentional blindness. For a conventional differentiable manifold, a second or higher order tangent space would give a better approximation to the local manifold structure than a simple plane (Pohl, 1962).

It is possible to introduce the evolutionary selection pressures of market forces into this model, using the approach of Wallace (2002).

Nonetheless, inattentional blindness, while constrained by multitasking, is not eliminated by it. This suggests that higher order institutional or machine cognition, the generalization of individual consciousness, is subject to canonical and idiosyncratic patterns of failure analogous to, but perhaps more subtle than, the kind of disorders described in Wallace (2005b, 2006a). Indeed, while machines designed along these principles – i.e. multitasking Global Workspace devices – could be spectacularly efficient at many complex tasks, ensuring their stability might be even more difficult than for institutions having the benefit of many centuries of cultural evolution.

In addition the necessity of interaction – synchronous or asynchronous – between internal giant components suggests the possibility of failures governed by the Rate Distortion Theorem. Forcing rapid communication between internal giant components ensures high error rates. Recent, and very elegant, ethnographic work by Cohen et al. (2006) and Laxmisan et al. (2006) regarding systematic medical error in emergency rooms focuses particularly on 'handover' problems at shift change, where incoming medical staff are rapidly briefed by outgoing staff. Systematic information overload in such circumstances seems almost routine, and is widely recognized as a potential error source within institutions.

This paper generalizes the Global Workspace model of individual consciousness to an analogous second order treatment of machine hyperconsciousness, and suggests, in particular, that multiple workspace multitasking significantly reduces, but cannot eliminate, the likelihood of inattentional blindness, of overfocus on one task to the exclusion of other powerful patterns of threat or affordance. It further appears

that rate distortion failure in communication between individual global workspaces will be potentially a serious problem for such systems - synchronous or sequential versions of the telephone game. Thus the multitasking hierarchical cognitive model appropriate to institutional or hyperconscious machine cognition is considerably more complicated than the equivalent for individual human consciousness, which seems biologically limited to a single shifting, tunable giant component structure. Human institutions, by contrast, appear able to entertain several, and perhaps many, such global workspaces simultaneously, although these generally operate at a much slower rate than is possible for individual consciousness. Hyperconscious machines, according to this model, would be able to function as efficiently as large institutions, but at or near the rate characteristic of individual consciousness.

Shared culture, however, seems to provide far more than merely a shared language for the establishment of the human organizations which enable our adaptation to, or alteration of, our varied environments. It also may provide the stabilizing mechanisms needed to overcome many of the canonical and idiosyncratic failure modes inherent to such structures - the embedding directives of law, tradition, and custom which have evolved over many centuries. Culture is truly as much a part of human biology as the enamel on our teeth (Richerson and Boyd, 2004). No such secondary heritage system is available for machine stabilization.

In sum, this paper contributes to a mathematical formalization of a machine hyperconsciousness based on a necessary conditions communication theory model quite similar to Dretske's attempts at understanding high level mental function for individuals. However, high order, multiple global workspace cognition seems not only far more complicated than is the case for individual animal consciousness, but appears prone to particular collective errors whose minimization, for institutions which operate along these principles, may have been the subject of a long period of cultural development.

Absent a detailed understanding of institutional stabilization mechanisms, and hence lacking a parallel engineering strategy for the construction of 'ethical' machines, reliable hyperconsciousness technology may simply not be possible.

MATHEMATICAL APPENDIX

The Shannon-McMillan Theorem

According to the structure of the underlying language of which a message is a particular expression, some messages are more 'meaningful' than others, that is, are in accord with the grammar and syntax of the language. The Shannon-McMillan or Asymptotic Equipartition Theorem, describes how messages themselves are to be classified.

Suppose a long sequence of symbols is chosen, using the output of the random variable X above, so that an output sequence of length n , with the form

$$x_n = (\alpha_0, \alpha_1, \dots, \alpha_{n-1})$$

has joint and conditional probabilities

$$P(X_0 = \alpha_0, X_1 = \alpha_1, \dots, X_{n-1} = \alpha_{n-1})$$

$$P(X_n = \alpha_n | X_0 = \alpha_0, \dots, X_{n-1} = \alpha_{n-1}).$$

(19)

Using these probabilities we may calculate the conditional uncertainty

$$H(X_n | X_0, X_1, \dots, X_{n-1}).$$

The uncertainty of the *information source*, $H[\mathbf{X}]$, is defined as

$$H[\mathbf{X}] \equiv \lim_{n \rightarrow \infty} H(X_n | X_0, X_1, \dots, X_{n-1}).$$

(20)

In general

$$H(X_n | X_0, X_1, \dots, X_{n-1}) \leq H(X_n).$$

Only if the random variables X_j are all stochastically independent does equality hold. If there is a maximum n such that, for all $m > 0$

$$H(X_{n+m} | X_0, \dots, X_{n+m-1}) = H(X_n | X_0, \dots, X_{n-1}),$$

then the source is said to be of *order* n . It is easy to show that

$$H[\mathbf{X}] = \lim_{n \rightarrow \infty} \frac{H(X_0, \dots, X_n)}{n+1}.$$

In general the outputs of the $X_j, j = 0, 1, \dots, n$ are *dependent*. That is, the output of the communication process at step n depends on previous steps. Such serial correlation, in fact, is the very structure which enables most of what is done in this paper.

Here, however, the processes are all assumed stationary in time, that is, the serial correlations do not change in time, and the system is *stationary*.

A very broad class of such self-correlated, stationary, information sources, the so-called *ergodic* sources for which the long-run relative frequency of a sequence converges stochastically to the probability assigned to it, have a particularly interesting property:

It is possible, in the limit of large n , to divide all sequences of outputs of an ergodic information source into two distinct sets, S_1 and S_2 , having, respectively, very high and very

low probabilities of occurrence, with the source uncertainty providing the splitting criterion. In particular the Shannon-McMillan Theorem states that, for a (long) sequence having n (serially correlated) elements, the number of ‘meaningful’ sequences, $N(n)$ – those belonging to set S_1 – will satisfy the relation

$$(21) \quad \frac{\log[N(n)]}{n} \approx H[\mathbf{X}].$$

More formally,

$$(22) \quad \begin{aligned} \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n} &= H[\mathbf{X}] \\ &= \lim_{n \rightarrow \infty} H(X_n | X_0, \dots, X_{n-1}) \\ &= \lim_{n \rightarrow \infty} \frac{H(X_0, \dots, X_n)}{n+1}. \end{aligned}$$

Using the internal structures of the information source permits *limiting attention only to high probability ‘meaningful’ sequences of symbols.*

The Rate Distortion Theorem

The Shannon-McMillan Theorem can be expressed as the ‘zero error limit’ of the Rate Distortion Theorem (Dembo and Zeitouni, 1998; Cover and Thomas, 1991), which defines a splitting criterion that identifies high probability pairs of sequences. We follow closely the treatment of Cover and Thomas (1991).

The origin of the problem is the question of representing one information source by a simpler one in such a way that the least information is lost. For example we might have a continuous variate between 0 and 100, and wish to represent it in terms of a small set of integers in a way that minimizes the inevitable distortion that process creates. Typically, for example, an analog audio signal will be replaced by a ‘digital’ one. The problem is to do this in a way which least distorts the *reconstructed* audio waveform.

Suppose the original stationary, ergodic information source Y with output from a particular alphabet generates sequences of the form

$$y^n = y_1, \dots, y_n.$$

These are ‘digitized,’ in some sense, producing a chain of ‘digitized values’

$$b^n = b_1, \dots, b_n,$$

where the b -alphabet is much more restricted than the y -alphabet.

b^n is, in turn, *deterministically retranslated* into a reproduction of the original signal y^n . That is, each b^n is mapped on to a unique n -length y -sequence in the alphabet of the information source Y :

$$b^n \rightarrow \hat{y}^n = \hat{y}_1, \dots, \hat{y}_n.$$

Note, however, that many y^n sequences may be mapped onto the *same* retranslation sequence \hat{y}^n , so that information will, in general, be lost.

The central problem is to explicitly minimize that loss.

The retranslation process defines a new stationary, ergodic information source, \hat{Y} .

The next step is to define a *distortion measure*, $d(y, \hat{y})$, which compares the original to the retranslated path. For example the *Hamming distortion* is

$$d(y, \hat{y}) = 1, y \neq \hat{y}$$

$$d(y, \hat{y}) = 0, y = \hat{y}.$$

(23)

For continuous variates the *Squared error distortion* is

$$d(y, \hat{y}) = (y - \hat{y})^2.$$

(24)

There are many possibilities.

The distortion between paths y^n and \hat{y}^n is defined as

$$d(y^n, \hat{y}^n) = \frac{1}{n} \sum_{j=1}^n d(y_j, \hat{y}_j).$$

(25)

Suppose that with each path y^n and b^n -path retranslation into the y -language and denoted y^n , there are associated individual, joint, and conditional probability distributions

$$p(y^n), p(\hat{y}^n), p(y^n|\hat{y}^n).$$

The *average distortion* is defined as

$$D = \sum_{y^n} p(y^n) d(y^n, \hat{y}^n).$$

(26)

It is possible, using the distributions given above, to define the information transmitted from the incoming Y to the outgoing \hat{Y} process in the usual manner, using the Shannon source uncertainty of the strings:

$$I(Y, \hat{Y}) \equiv H(Y) - H(Y|\hat{Y}) = H(Y) + H(\hat{Y}) - H(Y, \hat{Y}).$$

If there is no uncertainty in Y given the retranslation \hat{Y} , then no information is lost.

In general, this will not be true.

The *information rate distortion function* $R(D)$ for a source Y with a distortion measure $d(y, \hat{y})$ is defined as

$$R(D) = \min_{p(y, \hat{y}); \sum_{(y, \hat{y})} p(y) p(y|\hat{y}) d(y, \hat{y}) \leq D} I(Y, \hat{Y}).$$

(27)

The minimization is over all conditional distributions $p(y|\hat{y})$ for which the joint distribution $p(y, \hat{y}) = p(y)p(y|\hat{y})$ satisfies the average distortion constraint (i.e. average distortion $\leq D$).

The *Rate Distortion Theorem* states that $R(D)$ is the *maximum achievable rate of information transmission which does not exceed the distortion D* . Cover and Thomas (1991) or Dembo and Zeitouni (1998) provide details.

More to the point, however, is the following: Pairs of sequences (y^n, \hat{y}^n) can be defined as *distortion typical*; that is, for a given average distortion D , defined in terms of a particular measure, pairs of sequences can be divided into two sets, a high probability one containing a relatively small number of (matched) pairs with $d(y^n, \hat{y}^n) \leq D$, and a low probability one containing most pairs. As $n \rightarrow \infty$, the smaller set approaches unit probability, and, for those pairs,

$$p(y^n) \geq p(\hat{y}^n|y^n) \exp[-nI(Y, \hat{Y})].$$

(28)

Thus, roughly speaking, $I(Y, \hat{Y})$ embodies the splitting criterion between high and low probability pairs of paths.

For the theory of interacting information sources, then, $I(Y, \hat{Y})$ can play the role of H in the dynamic treatment that follows.

The rate distortion function can actually be calculated in many cases by using a Lagrange multiplier method – see Section 13.7 of Cover and Thomas (1991).

References

- Aiello W., F. Chung, and L. Lu, 2000, A random graph model for massive graphs, in *Proceedings of the 32nd Annual ACM Symposium on the Theory of Computing*.
- Albert R., and A. Barabasi, 2002, Statistical mechanics of complex networks, *Reviews of Modern Physics*, 74:47-97.
- Ash R., 1990, *Information Theory*, Dover Publications, New York.
- Atlan H., and I. Cohen, 1998, Immune information ,self-organization and meaning, *International Immunology*, 10:711-717.
- Auslander L., 1967, *Differential Geometry*, Harper and Row, New York.
- Baars B., 1988, *A Cognitive Theory of Consciousness*, Cambridge University Press, New York.
- Baars B., and S. Franklin, 2003, How conscious experience and working memory interact, *Trends in Cognitive Science*, doi:10.1016/S1364-6613(03)00056-1.
- Baars, B., 2005, Global workspace theory of consciousness: toward a cognitive neuroscience of human experience, *Progress in Brain Research*, 150:45-53.
- Bak A., R. Brown, G. Minian and T. Porter, 2006, Global actions, groupoid atlases and related topics, *Journal of Homotopy and Related Structures*, 1:1-54. Available from ArXiv depository.
- Bennett M., and P. Hacker, 2003 *Philosophical Foundations of Neuroscience*, Blackwell Publishing, London.
- Burago D., Y. Burago, and S. Ivanov, 2001, *A Course in Metric Geometry*, American Mathematical Society, Providence, RI.
- Cohen I., 2000, *Tending Adam's Garden: Evolving the Cognitive Immune Self*, Academic Press, New York.
- Cohen T., B. Blatter, C. Almeida, E. Shortliffe, and V. Patel, 2006, A cognitive blueprint of collaboration in context: Distributed cognition in the psychiatric emergency department, *Artificial Intelligence in Medicine*, 37:73-83.
- Connes A., 1994, *Noncommutative Geometry*, Academic Press, San Diego.
- Corless R., G. Gonnet, D. Hare, D. Jeffrey, and D. Knuth, 1996, On the Lambert W function, *Advances in Computational Mathematics*, 4:329-359.
- Cover T., and J. Thomas, 1991, *Elements of Information Theory*, John Wiley and Sons, New York.
- Dehaene S., and L. Naccache, 2001, Towards a cognitive neuroscience of consciousness: basic evidence and a workspace framework, *Cognition*, 79:1-37.

- Dehaene S., and J. Changeux, 2005, Ongoing spontaneous activity controls access to consciousness: a neuronal model for inattentive blindness, *PLOS Biology*, 3:e141.
- Dembo A., and O. Zeitouni, 1998, *Large Deviations: Techniques and Applications*, Second edition, Springer, New York.
- Dretske F., 1981, *Knowledge and the Flow of Information*, MIT Press, Cambridge, MA.
- Dretske F., 1988, *Explaining Behavior*, MIT Press, Cambridge, MA.
- Dretske, F., 1993, Mental events as structuring causes of behavior, in *Mental Causation* (ed. by A. Mele and J. Heil), pp. 121-136, Oxford University Press.
- Dretske F., 1994, The explanatory role of information, *Philosophical Transactions of the Royal Society A*, 349:59-70.
- Erdos P., and A. Renyi, 1960, On the evolution of random graphs, reprinted in *The Art of Counting*, 1973, 574-618 and in *Selected Papers of Alfred Renyi*, 1976, 482-525.
- Feynman, R., 1996, *Feynman Lectures on Computation*, Addison-Wesley, Reading, MA.
- Freeman, W., 2003, The wave packet: an action potential of the 21st Century, *Journal of Integrative Neurosciences*, 2:3-30.
- Fullilove, M., 2004, *Root Shock*, Ballantine Books, New York.
- Glazebrook, J., 2006, Rate distortion manifolds as model spaces for cognitive information. In preparation.
- Golubitsky M., and I. Stewart, 2006, Nonlinear dynamics and networks: the groupoid formalism, *Bulletin of the American Mathematical Society*, 43:305-364.
- Goubault, E., and M. Raussen, 2002, Dihomotopy as a tool in state space analysis, *Lecture Notes in Computer Science*, Vol. 2286, April, 2002, pp. 16-37, Springer, New York.
- Goubault E., 2003, Some geometric perspectives in concurrency theory, *Homology, Homotopy, and Applications*, 5:95-136.
- Granovetter M., 1973, The strength of weak ties, *American Journal of Sociology*, 78:1360-1380.
- Grimmett G., and A. Stacey, 1998, Critical probabilities for site and bond percolation models, *The Annals of Probability*, 4:1788-1812.
- Gunderson L., 2000, Ecological resilience - in theory and application, *Annual Reviews of Ecological Systematics*, 31:425-439.
- Holling C., 1973, Resilience and stability of ecological systems, *Annual Reviews of Ecological Systematics*, 4:1-23.
- Laxmisan A., F. Hakimzada, O. Sayan, R. Green, J. Zhang, V. Patel, 2006, The multitasking clinician; Decision-making and cognitive demand during and after team handoffs in emergency care. Submitted.
- Khinchin A., 1957, *The Mathematical Foundations of Information Theory*, Dover Publications, New York.
- Kozma R., M. Puljic, P. Balister, B. Bollobas, and W. Freeman, 2004, Neuroperturbation: a random cellular automata approach to spatio-temporal neurodynamics, *Lecture Notes in Computer Science*, 3305:435-443.
- Kozma R., M. Puljic, P. Balister, and B. Bollobas, 2005, Phase transitions in the neuroperturbation model of neural populations with mixed local and non-local interactions, *Biological Cybernetics*, 92:367-379.
- Krebs, P., 2005, Models of cognition: neurological possibility does not indicate neurological plausibility, in Bara, B., L. Barsalou, and M. Bucciarelli (eds.), *Proceedings of CogSci 2005*, pp. 1184-1189, Stresa, Italy. Available at <http://cogprints.org/4498/>.
- Lee J., 2000, *Introduction to Topological Manifolds*, Springer, New York.
- Luczak T., 1990, *Random Structures and Algorithms*, 1:287.
- Mack A., 1998, *Inattentive Blindness*, MIT Press, Cambridge, MA.
- Matsuda, T., and R. Nisbett, 2006, Culture and change blindness, *Cognitive Science*, 30:381-399.
- Molloy M., and B. Reed, 1995, A critical point for random graphs with a given degree sequence, *Random Structures and Algorithms*, 6:161-179.
- Molloy M., and B. Reed, 1998, The size of the giant component of a random graph with a given degree sequence, *Combinatorics, Probability, and Computing*, 7:295-305.
- Newman M., S. Strogatz, and D. Watts, 2001, Random graphs with arbitrary degree distributions and their applications, *Physical Review E*, 64:026118, 1-17.
- Newman M., 2003, Properties of highly clustered networks, arXiv:cond-mat/0303183v1.
- Patel, V., 1998, Individual to collaborative cognition: a paradigm shift? *Artificial Intelligence in Medicine*, 12:93-96.
- Pielou, E., 1977, *Mathematical Ecology*, John Wiley and Sons, New York.
- Pohl W., 1962, Differential geometry of higher order, *Topology* 1:169-211.
- Pratt V., 1991, Modeling concurrency with geometry, *Proceedings of the 18th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 311-322.
- Richerson P., and R. Boyd, 2004, *Not by Genes Alone: How Culture Transformed Human Evolution*, Chicago University Press.
- Savante J., D. Knuth, T. Luczak, and B. Pittel, 1993, The birth of the giant component, arXiv:math.PR/9310236v1.
- Shannon C., and W. Weaver, 1949, *The Mathematical Theory of Communication*, University of Illinois Press, Chicago, IL.
- Simons, D., and C. Chabris, 1999, Gorillas in our midst: sustained inattentive blindness for dynamic events, *Perception*, 28:1059-1074.
- Simons D., 2000, Attentional capture and inattentive blindness, *Trends in Cognitive Sciences*, 4:147-155.
- Stewart I., M. Golubitsky, and M. Pivato, 2003, Symmetry groupoids and patterns of synchrony in coupled cell networks, *SIAM Journal of Applied Dynamical Systems*, 2:609-646.
- Stewart I., 2004, Networking opportunity, *Nature*, 427:601-604.
- Tononi G., 2004, An information integration theory of consciousness, *BMC Neuroscience*, 5:42.
- Wallace R., 2000, Language and coherent neural amplification in hierarchical systems: renormalization and the dual information source of a generalized spatiotemporal stochastic

resonance, *International Journal of Bifurcation and Chaos*, 10:493-502.

Wallace R., 2002, Adaptation, punctuation, and information: a rate distortion approach to non-cognitive ‘learning plateaus’ in evolutionary process, *Acta Biotheoretica*, 50:101-116.

Wallace R., 2005a, *Consciousness: A Mathematical Treatment of the Global Neuronal Workspace Model*, Springer, New York.

Wallace R., 2005b, A global workspace perspective on mental disorders, *Theoretical Biology and Medical Modelling*, 2:49, <http://www.tbiomed.com/content/2/1/49>.

Wallace R., 2005c, The sleep cycle: a mathematical analysis from a global workspace perspective, <http://cogprints.org/4517/>

Wallace R., 2006a, Pitfalls in biological computing: canonical and idiosyncratic dysfunction of conscious machines, *Mind and Matter*, 4:91-113.

Wallace R., 2006b, Institutional cognition, <http://cogprints.org/xxxx/>.

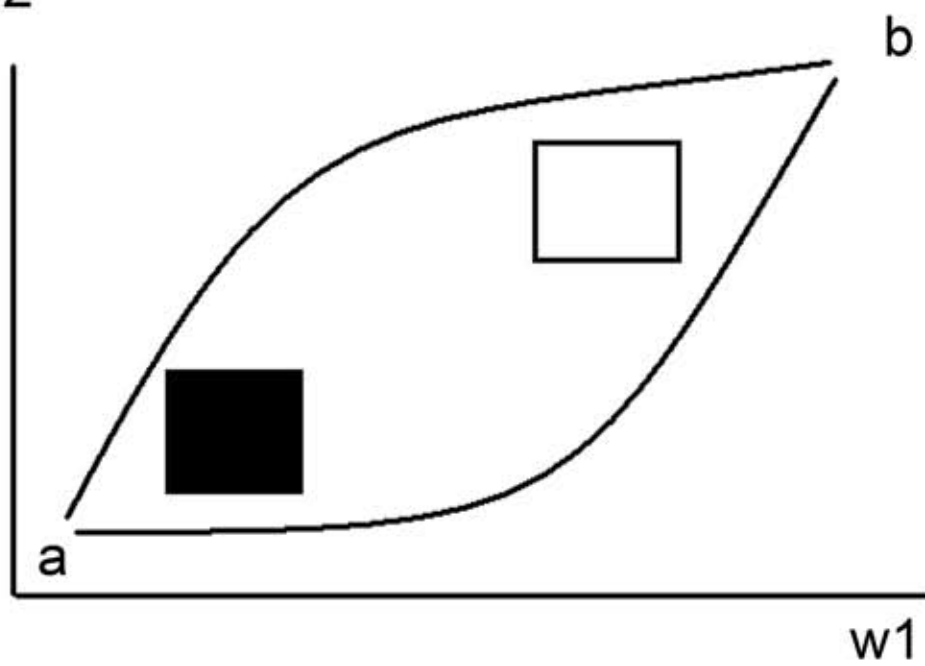
Wayand J., D. Levin and D. Alexander Varakin, 2005, Inattentional blindness for a noxious multimodal stimulus, *American Journal of Psychology*, 118:339-352.

Weinstein A., 1996, Groupoids: unifying internal and external symmetry, *Notices of the American Mathematical Association*, 43:744-752.

Figure Captions

Figure 1a. Diagonal Black and White holes in the two dimensional ω -plane. Only two direct paths can link points a and b which are continuously deformable into one another without crossing either hole. There are two additional monotonic switchback paths which are not drawn.

Figure 1b. Cross-diagonal Black and White holes as in 1a. Three direct equivalence classes of continuously deformable paths can link a and b . Thus the two spaces are topologically distinct. Here monotonic switchbacks are not possible, although relaxation of that condition can lead to ‘backwards’ switchbacks and intermediate loopings.

w_2  w_2 