

# ‘The frozen accident’ as an evolutionary adaptation: A rate distortion theory perspective on the dynamics and symmetries of genetic coding mechanisms

–Version 2

James F. Glazebrook\*, PhD

Department of Mathematics and Computer Science

Eastern Illinois University

and

Rodrick Wallace†, PhD

Division of Epidemiology

The New York State Psychiatric Institute

February 21, 2011

## Abstract

We survey some interpretations and related issues concerning the frozen hypothesis due to F. Crick and how it can be explained in terms of several natural mechanisms involving error-correction codes, spin glasses, symmetry breaking and the characteristic robustness of genetic networks. The approach to most of these questions involves using elements of Shannon’s rate distortion theory incorporating a semantic system which is meaningful for the relevant alphabets and vocabulary implemented in transmission of the genetic code. We apply the fundamental homology between information source uncertainty with the free energy density of a thermodynamical system with respect to transcriptional regulators and the communication channels of sequence/structure in proteins. This leads to the suggestion that the frozen accident may have been a type of evolutionary adaptation.

**Key words** Frozen accident, rate distortion function, protein folding, free energy density, spin glass, groupoid, Onsager relations, holonomy.

## 1 Introduction

Examining and predicting the geometric/topological structures of the genetic coding network is essential to understanding its (co)evolution as a complex communications system, employing a vocabulary of a given genetic code that determines the family of proteins encodable by the genes themselves. The architecture of this network developed from a coevolution of genes and of genetic structures that were progressively conditioned to shield against translation and replication errors. Crick’s hypothesis (Crick, 1966, 1968; surveyed in e.g. Ardell and Sella, 2002), in broad terms, says that on reading the mRNA, the code determines the amino acid sequence of the evolved proteins, as is the case for most organisms, so in a post-transitional phase any kind of alteration to the size of the code would have dire consequences owing to a global impact on proteins created by new amino acids subject to possible nonsensical messaging. Crick gave flexible rules for pairing the third base of the codon with the first base of the anticodon, to the extent that a single tRNA type would be able to recognize up to three codons. More complex protein structures arise when there is an enrichment and expansion of the vocabulary while any ambiguity in the code is minimized, so restricting the content of information. When the codon meaning is altered, the information selected would condition that codon to some advantage. In this way the “freezing” was professed to be an outcome of such selective restrictions and this would pull the reins on further evolvability. While over the years there has been much debate and challenge concerning these rules and to establish a concrete mechanism for the companion “wobble

---

\*Department of Mathematics and Computer Science, Eastern Illinois University, 600 Lincoln Avenue, Charleston IL 61920–3099 USA, email [jfglazebrook@eiu.edu](mailto:jfglazebrook@eiu.edu)

†549 W. 123 St., Suite 16F, New York, NY, 10027, USA. Telephone (212) 865-4766, email [rdwall@ix.netcom.com](mailto:rdwall@ix.netcom.com). Affiliation is for identification only.

hypothesis”, we outline here several scenarios from the point of view of coevolutionary rate distortion dynamics in graphs that can exhibit ‘robustness’ while admitting ‘meaningful’ signalling paths which are susceptible to vocabulary enrichment, and further, admit structure preserving patterns that evolve towards optimizing error-correction. These collective mechanisms can be formulated in the context of a spin-glass model (cf Ciliberti et al., 2007), that incorporates the Onsager relations of statistical physics applied to networks of mutating sequences and error-correction in the presence of rate distortion dynamics, then leading to phase transitions through which symmetry breaking occurs and results in changing the graph’s topological structure. These observations are supported by a number of theoretical findings which, in fact, are relatively recent and so it befits us to provide here some of the necessary background material.

Firstly, a position often maintained is that evolution influences the emergence of the genetic code by selecting an amino acid map that is error-minimizing and subsequent competition between organisms is determined by the overall capability of their respective codes. Following this line of thought, Tlusty (2007, 2008a, 2008b), implementing a topological graph-theoretic approach, has developed a model for the emergence of the genetic code as a supercritical phase transition occurring within noisy information channels as traced by maps between nucleotides and amino acids with error bounds in place. The proposed paradigm is that these processes are indeed ‘cognitive’ (Wallace, 2010a, 2010b, 2010c; Wallace and Wallace, 2009) following the immunology/language perspective of Atlan and Cohen (1998) (see also, Cohen 2000; Cohen and Harel, 2007) that human and biological organizations at all scales are cognitive in so far that patterns of threat and opportunity are perceived, these patterns are compared with an internal image of the environment, and then a choice of responses from a vast repertoire of possibilities is initiated.

This present paper continues with this theme to establish one of several possible corollaries derived from Wallace (2010a, 2010c) by addressing the question of coevolutionary robustness against errors, error-correction and phase transitions modeled by the topological dynamics of graphs that can represent certain spin glass/error correcting structures that are susceptible to thermodynamic spontaneous symmetry breaking, thus shedding further light on the question of what exactly was the “accident” that did occur. Such symmetry breaking of the genetic code, has been considered in the context of Lie algebra representations by Bashford et al. (1998, 2006) (cf Hornos and Hornos, 1994). One perspective taken up here, using rate distortion dynamics, is that such a sequence of broken symmetries corresponds to phase transitions in the underlying error correcting networks through which the codon allocation to amino acids is mainly the outcome of error-correction minimization and efficiency (see Bashford et al., 1998, and relevant references therein), a scenario that appears to be on par with the approach Ardell and Sella (2002), Sella and Ardell (2002, 2006).

While on the mathematical-physical side of things, several explanations for “freezing” and “wobbling” can be given in terms of error-correction and the structural theory of Lie algebras, which we survey, a novel technique introduced here involves showing how the dynamics governing the underlying mechanisms can be represented in terms of a ‘covariant differentiation’ of the Shannon entropy along ‘meaningful paths’ embedded in a genetic graph which may have, among other things, a correlation with error-correction and folding rates. This operation over which the various ‘directions’ are taken \* subsequently determines the *holonomy* of the system through an error-correction network—a broader scale geometric representation of transitional phases in which the broken symmetries may be expressed in terms of holonomy groups that collectively, via disjoint union, form a holonomy groupoid which in principle can be made explicit.

## 2 ‘The frozen accident’— or not quite

We start by putting matters into perspective by surveying some basic observations. Recall that genes can be represented by molecular words written in terms of the nucleotide bases  $U, C, G$  and  $A$ , whereas proteins are written in a language of 20 letters corresponding to the amino acids in which each of the latter is encoded by specific triplets of the basis members, known as *codons*. In theory there are  $64 = 4^3$  codons with the number of possible observables lying somewhere between 48 and 64 (see e.g. Koonin and Novozhilov, 2009; Tlusty 2007); however, it is claimed in Koonin and Novozhilov (2009) that the code mapping the 64 codons to the 20 amino acids is anything but random. There are at least 48 discernable codons but only 20 amino acids available (and 3 stop codons), so the code is degenerate in so far that several codons can represent the same amino acid. Entropy analysis (Adami et al., 2000; Lisewski, 2008) reveals that the information content of a random protein structure can occupy  $\log_2(20) \simeq 4.32$  bits of entropy per amino acid residue in a primary sequence.

In the presence of topological changes there would have been alterations of an excessive amount of (protein) structures and those frequently observed tend to be the ones that have managed to remain intact as the structures became more complex. The “wobble rules” assume that only 48 codons can be distinguished owing to the physiochemical limitations of the translational mechanism and the resulting codon graph converges to 20 amino acids. The question is: does a single sRNA molecule recognize several codons? The “wobble” effect aside, there exist 64 distinguishable codons and the maximal number of amino acids increases to 25 which is not a dramatic amount by any means, though it has been a puzzling matter as to

---

\*A reader with some acquaintance with differential geometry will understand this as ‘covariant differentiation over (or along) a vector field’— an operation specified by choice of ‘connection’. This we implement on graphs in §6.

why evolution did freeze prior to improving the translational mechanism to single out all 64 codons. Once the meaning of a codon had changed, again, selectivity would apply that codon to a site for a new amino acid to serve to some advantage, or otherwise simply to replace it.

The traditional approach to producing more tRNAs would have been to change the anticodons of existing ones, giving rise to a new class of amino acid proliferating across the code while systematically reshuffling a large number of codons in the process. In more basic terms, interfering with the genetic code could change the meaning of a codon, thus reducing, from our viewpoint, the fidelity of information when the rate distortion estimate is violated (see §3.1). As was recalled in the introduction, Crick's hypothesis had suggested that no new amino acids could arise without disrupting a large number of proteins, hence stalling evolution – a claim that has since been challenged from many fronts (see e.g. Ardell and Sella, 2002; Söll and RajBhandary, 2006). A product of the coevolutionary dynamics gives rules for load minimization and diversification for regulating patterns of the code that were robust to both error and redundancy, the degrees of which are influenced by the code's topology that would have been alterable through sequences of stochastic fluctuations. Codons interchanged through error may subsequently be assigned to compatible amino acids so minimizing the possible detrimental effects. At the same time, an enrichment of the vocabulary provided a broader scope for the encoding of proteins (see Sella and Ardell, 2006).

Vetsigian et al. (2006) claim a 'communality' and 'universality' to be established out of a tournament between a variety of innovative sharing protocols which may include several non-Darwinian mechanisms. Relative to time scales, the long-term reduces ambiguity, whereas in the short-term the code has to be fortified to tolerate a higher degree of ambiguity in assimilating new types of genes. More specifically, from Vetsigian et al. (2006):

A protein that is robust to translational errors *a fortiori* is also more tolerant to translation with a different code. Conversely, the less optimized the recipient code, the more error-tolerant its proteins, and therefore the less harmful the effect on the established genes of a code change in the direction of the donor code. This has the important consequence that in the initial stages of the genetic code evolution, when the diversification tendency of codes was strongest, HGT (horizontal gene transfer) was possible and must have been extensive despite the presence of many different codes ... Once the optimization of the genetic code is complete, there is no pressure to maintain compatibility. Therefore, the "freezing" of the universal genetic code could trigger the radiation of the underlying translational machineries...

We may reasonably assume that transmission errors eventually corrupt code patterns and those codes that can withstand and manipulate errors possess natural advantages over those that do not. Ardell and Sella (2002) in concluding differently to Crick's assertion, perceive code-message evolution as producing structure preserving codes which have near optimal error-correcting properties, with the selection of mutations and translational error inducing a bias in the codon distribution to amino acids which in the long-term favors optimal error-correction patterns. Crick's claim of "freezing" makes some sense because the errors themselves condition evolution to some sort of frozen state of an error-correcting code. Specifically, the claim is that an evolutionary constraint on messages with respect to selective pressures, may actually induce the error-correcting codes to evolve rather than to have erased them altogether. Thus, in this evolutionary context the allied and relevant mechanisms of protein synthesis, folding and mutations, provide suitable clues.

An underlying assumption proposed in Adami et al. (2000) is that an organism's complexity reflects upon that of its genome and hence has evolutionary consequences. So one may ask what actually is the information provided by DNA beyond a road map for the structure of an organism? The current perspective sees this as a blueprint for constructing an organism that can survive within its native environment and then pass on that information to its progeny (cf Dawkins, 1976). In this respect, an organism's DNA catalogs not only information concerning its structure, but to some extent information concerning its environment and the coevolution of its species as well. In keeping with this basic principle, one may propose an explanation of genomic complexity within the framework of Shannon's basic principles (see Adami et al., 2000; Adami, 2004, and references therein for related work). It is in this respect that the fundamental theorems of information transmission are sufficiently general to the extent that biological systems can sustain a Shannon-based coding scheme to facilitate the transmission of genomic information within a range of mechanisms, provided that semantics can to be incorporated as a functional component (see §4.1 and cf Dretske, 1981).

## 3 Encoding and Decoding

### 3.1 The rate distortion function

For the sake of self-containment in this paper, we briefly recall some elementary facts from the Shannon theory. As it is commonly understood, distortion arises when there is a fast relay of information through some channel which exceeds the latter's capacity. One of the guiding principles asserts that in order to reproduce a message transmitted from a source to a receiver, it is necessary to know what sort of information should be transmitted and how. These facts along with

specifying the nature of the communicating channel are essential ingredients for engineering a reliable encoding/decoding system. Following Berger (1971) we briefly recall some of the basic operations.

**Source encoder:** We may consider some output  $x(t)$  emanating from the source as projected to a finite set of preselected images; namely, the space of possible source outputs is partitioned into a set of *equivalence classes* and the source encoder informs the channel encoder of that class containing the particular source output observed. Once the channel encoder is informed that the source output belongs to say, the  $m$ -th equivalence class, it transforms the corresponding waveform  $\tilde{x}_m(t)$  across the channel. These equivalence classes are schematically represented by a graph or network are manifestly the main computational procedures as shown in this paper.

**Source decoder:** Within the system is a cascade of a channel encoder and a source decoder. The channel decoder receives a waveform  $\tilde{y}(t)$  of a corresponding function  $y(t)$  over some time interval and decides upon the nature of the message as transmitted. Then it sends its approximation  $m'$  of the message number to the source decoder which in turn creates  $y_{m'}(t)$  to register the system's estimate of  $x(t)$  over that time interval. Initially, we may think of  $x(t)$  and  $y(t)$  as 'waveforms', but in our case, we consider these as consisting of a language with its own intrinsic grammar/syntax, as well as 'meaning' – to be made more specific in §4.1. Analogous considerations apply to the channel signals  $\tilde{x}(t)$  and  $\tilde{y}(t)$ .

One of Shannon's notable results was that a communication system can be designed such that it achieves a level of fidelity  $D$  once the *rate distortion*  $R(D) \leq C$ , where  $C$  denotes the channel capacity. Putting it another way, if the receiver can tolerate an average amount of distortion  $D$ , the rate distortion  $R(D)$  is the effective rate at which the source can relay information with that level of tolerance, and the estimate  $R(D) \leq C$  is a necessary condition for effective communication. More specifically,  $R(D)$  can be defined in terms of *average mutual information* as follows. Firstly, for  $k, j$  running over a suitable alphabet, let us write a given conditional probability assignment as  $Q(k|j)$  such that in the usual way, we have an associated joint distribution  $P(j, k) = P(j)Q(k|j)$ . We express *the average distortion* as

$$d(Q) = \sum_{j,k} P(j)Q(k|j) d(j, k), \quad (3.1)$$

where  $d(, )$  denotes the distortion measure. A conditional probability assignment  $Q(k|j)$  is said to be *D-admissible* if and only if  $d(Q) \leq D$ . The set of all *D*-admissible conditional probability assignments we denote by

$$Q_D = \{Q(k|j) : d(Q) \leq D\}. \quad (3.2)$$

Along with an average distortion  $d(Q)$ , we also have an *average mutual information*

$$I(Q) = \sum_{j,k} P(j)Q(k|j) \log \left[ \frac{Q(k|j)}{Q(k)} \right]. \quad (3.3)$$

Then for fixed  $D$ , the rate distortion function is defined as

$$R(D) = \min_{Q \in Q_D} I(Q). \quad (3.4)$$

The rate at which a source produces information subject to insisting upon perfect reproduction, is the *source entropy*  $H$ . Given a distortion measure such that perfect reproduction is assigned zero distortion, then we have  $R(0) = H$ . As  $D$  increases,  $R(D)$  becomes a monotonically decreasing (convex) function which eventually is zero, typically at a maximum value for  $D$  (see Berger, 1971, Ch. 1). This is a very basic observation, and typically in rate distortion theory one seeks a reduction of  $H$  by either slowing down the emission of coding, or encoding the relevant languages at a lower rate. In view of Shannon's theorem, as long as  $H < C$ , there will be suitable fidelity in transmission. For the case of genetic coding considered here, conditions of *discrete memoryless information source* (DMI) and *discrete memoryless channels* (DMC) (Yockey, 2005; McEliece, 2004) are usually assumed, but in any event, how well a communicating system can evolve in order to satisfy such an estimate is a common problem for communications engineering since in practice the source rate may be corrupted due to low memory and coding congestion; for protein folding and mutations; references Adami, (2004), Crooks and Brenner (2004), Lisewski (2008), Thustly (2007, 2008), Wallace (2010a, 2010c) address such issues.

### 3.2 The Groupoid Free Energy Density

Recall that for a thermodynamic state of a given system at fixed temperature  $T$  with energy  $E$  and entropy  $S$ , the *free energy density*  $F$  is defined to be

$$F = E - TS. \quad (3.5)$$

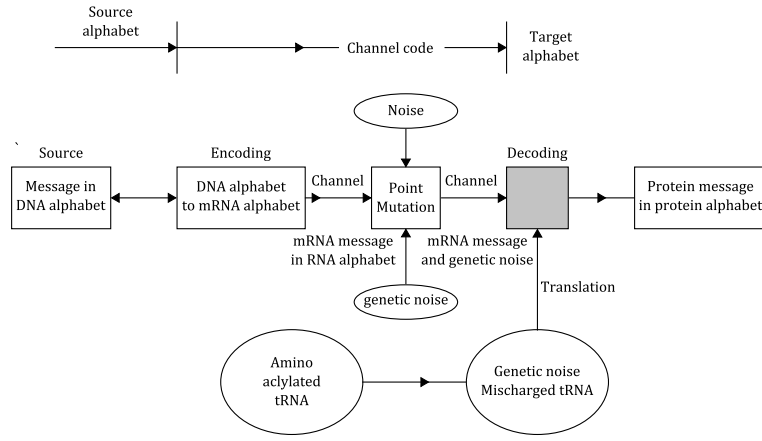


Figure 1: Outline of transmission of genetic messages from the DNA tape to the protein tape as based on Fig. 5.2 (p.31) of Yockey (2005).

In the Hamiltonian formalism one takes the volume  $V$  and the partition function  $Z(K)$  derived from the system's Hamiltonian at inverse temperature  $K$  (Landau and Lifshitz, 2007; Kurzynski, 2006). The free energy density is then defined to be

$$\begin{aligned}
 F[K] &= \lim_{V \rightarrow \infty} -\frac{1}{K} \frac{\log[Z(K, V)]}{V} \\
 &= \lim_{V \rightarrow \infty} \frac{\log[\hat{Z}(K, V)]}{V}, \text{ where } \hat{Z} = Z^{-\frac{1}{K}}.
 \end{aligned} \tag{3.6}$$

At this stage we introduce the *groupoid* concept (generalizing the algebraic concept of a 'group') in relationship to *equivalence classes* (which can be based upon a network with concatenation of edges), as explained in Appendix §A.1 (see also Glazebrook and Wallace, 2009a, 2009b). Thus consider an information source  $H_{G_\alpha}$  over a corresponding groupoid  $G_\alpha$ ; heuristically, we can consider  $H$  as parametrized by  $G_\alpha$ . The probability of  $H_{G_\alpha}$  is given by:

$$P(H_{G_\alpha}) = \frac{\exp[-H_{G_\alpha} K]}{\sum_{\beta} \exp[-H_{G_\beta} K]}, \tag{3.7}$$

where the normalizing sum is over all possible subgroupoids of the largest available symmetry groupoid. On setting

$$Z_G = \sum_{\alpha} \exp[-H_{G_\alpha}], \tag{3.8}$$

the *groupoid free energy density (GFE)* of the system  $F_G$  at inverse normalized equivalent temperature  $K$  is then defined as

$$F_G[K] = -\frac{1}{K} \log[Z_G(K)]. \tag{3.9}$$

With each such groupoid  $G_\alpha$  we can associate a dual information source  $H_{G_\alpha}$ . We recall the rate distortion function between the message sent by the cognitive process and the observed impact, while noting that both  $H_{G_\alpha}$  and  $R(D)$  may be considered as free energy density measures. In a sense,  $R(D)$  constitutes a sort of 'thermal bath' for the process of cognition. Then the probability of the dual information source can be expressed by

$$P(H_{G_\alpha}) = \frac{\exp[-H_{G_\alpha}/\kappa R(D)\tau]}{\sum_{\beta} \exp[-H_{G_\beta}/\kappa R(D)\tau]}, \tag{3.10}$$

where  $\kappa$  denotes a suitable dimensionless constant characteristic of the system in the context of a fixed 'machine response time'  $\tau$ . Associated with (3.10) is a *free energy Morse Function*

$$F_R = -\lambda R(D) \log\left[\sum_{\alpha=1}^n \exp[-H_\alpha/\lambda R(D)]\right], \tag{3.11}$$

whose critical point behavior determines certain topological characteristics of an underlying manifold that can be expressed in terms of its Morse-theoretic indices (Matsumoto, 2001). In each case the sum is over all possible subgroupoids of the

largest available symmetry groupoid (see Appendix §A.1). Accordingly, the term  $R(D)\kappa$  in (3.10) represents a rate distortion energy, in this case, a kind of temperature analog. In the context of a fixed  $\tau$ , a decline in  $R(D)$  (on increase in average distortion), acts to ‘lower the machine temperature’ and thus driving it to more simple albeit less enriched signalling. Observe that if a range over all possible  $\alpha$  is taken, the groupoids  $G_\alpha$  and corresponding relationships such as (3.10), create an even larger picture which involves the structure of a *groupoid atlas* (see Bak et al., 2006) which has already been applied to several descriptive cognitive mechanisms as seen in Glazebrook and Wallace (2009a, 2009b, 2010).

### 3.3 Phase transition and symmetry breaking

The relation between phase transitions in physical systems and topological changes has become a central topic of research across a broad range of subdisciplines. One can see that phase transitions in physical systems are ubiquitous, following Landau’s group symmetry shifting arguments (Landau and Lifshitz, 2007; Pettini, 2007). Higher temperatures enable higher system symmetries, and as temperature changes, punctuated shifts to different symmetry states occur in characteristic manners. Franzosi and Pettini (2004) argue that the standard way of studying phase transitions in physical systems is to consider how the empirical values of thermodynamic states, vary with temperature, volume, or an external field, and then to associate the experimentally observed discontinuities at a phase transition to the occurrence of a singularity. At such an instance analyticity fails in the mathematical sense though it remains to be seen whether this is the ultimate level of an analytic understanding of such transitional phenomena, or if indeed some reduction to a more basic level is possible. It is observed that non-analyticity is the “shadow” of a more fundamental phenomenon occurring in a given model space: *a topology change*, and that the latter is a *necessary* condition for a phase transition to occur. Such topology changes can be studied within the framework of Morse theory influenced topological structures such as the case, say, for certain handle-body decompositions (Matsumoto, 2001), an essential observation that may be consequential for protein functions (cf Wallace, 2010c). Note however, that the converse of the main result of Franzosi and Pettini (2004) does not hold, thus ruling out a one-to-one correspondence between phase transitions and topology changes. An open problem is that of *sufficiency* conditions, that is, to determine which kinds of topology changes can influence a phase transition (and how). There are other approaches such as demonstrated in relatively straightforward models, where as in Rose et al. (1990) a fuzzy clustering system based of annealing through a probabilistic process leads to phase transitions with critical (non-zero) vectors for the free energy at each temperature.

Extension of such transitional arguments in terms of rate distortion and metabolic measures appear direct, particularly in the setting of the groupoids constructed by the disjoint union of the homology groups representing the different coding topologies identified in Tlustý (2007a) (see Wallace, 2010a). To clarify matters, let us recall that in many thermodynamic systems, the associated Hamiltonian may be invariant under a symmetry transformation due to certain parameter changes, in contrast to the lowest energy state which is not. In subsequent phase transitions the overall symmetry is lost (*spontaneous symmetry breaking*) and consequently lower temperature states will admit lower symmetries, and due to the randomization of higher temperatures, the higher states will become more accessible to the system as a result of their modified symmetries and energy levels (Landau and Lifshitz, 2007). In the informational context of error-correction, we will need to turn to the fundamental homology between the Shannon entropy and the free energy density of the system as outlined in §4.1.

This scenario becomes more apparent when we look at the symmetries of the genetic code. For instance, Hornos and Hornos (1994) recall from Bertman and Jungck (1978) the computation of at least  $10^{71} - 10^{84}$  possible genetic codes by permuting the 64 codons and distributing them over 20 amino acids. By considering those Lie algebras admitting 64 dimensional irreducible representations, Hornos and Hornos (1994) (cf Bashford et al., 1998, 2008) initiate a chain of sub-representations commencing from the Lie algebra  $\mathfrak{sp}(6)$ , and postulate a sequence of symmetry breaking in accordance with that chain:

$$\mathfrak{sp}(6) \supset \mathfrak{sp}(4) \oplus \mathfrak{su}(2) \supset \mathfrak{su}(2) \oplus \mathfrak{su}(2) \oplus \mathfrak{su}(2) \supset \mathfrak{su}(2) \oplus \mathfrak{u}(1) \oplus \mathfrak{su}(2) \supset \mathfrak{su}(2) \oplus \mathfrak{u}(1) \oplus \mathfrak{u}(1). \quad (3.12)$$

At any stage the number of representations occurring corresponds to the number of amino acids that were then incorporated into the code and those currently observed are the net outcome of broken symmetries. In this analysis four amino acids (phenylalanine, serine, arginine and cysteine) seemingly do not divide under the  $U(1)$ -action. If they had subdivided they would have created a “symmetry perfect code” with 26 amino acids (hence a redundancy of 6) and a stop code (see Fig. 1 of Hornos and Hornos, 1994). Such a claim may be compared with the combinatorial-geometric arguments based on the topology of codon space in Tlustý (2007a) (see also §6.1) suggesting that further evolutionary measures may expand the code’s expression from 20 to possibly 25 amino acids.

The observations of Bashford et al. (1998, 2008) reflect back upon an earlier claim of Jukes (1983) that the “freezing” of the code would have been the result of partial symmetry breaking achieved by the aforementioned parameter choices in the Hamiltonian. Bashford et al. (1998, 2008) differ in their approach by opting for codon-anticodon pairings in place of codon-amino acid assignments and apply combinatorial-branching techniques commencing from the Lie algebra  $\mathfrak{sl}(6, 1)$ . Besides identifying possible “wobble-effects” due to reshuffling through combinatorial symmetries, they investigate the structure of eukaryotic and vertebrate mitochondrial codes along branching chains and introduce a  $\mathbb{Z}_2$ -grading on codon space (just as

there is a grading into bosonic and fermionic types in quantum statistics) thus extending matters towards representations of super Lie algebras. Along with these codes are variants such as metabacteria and chloroplast codes with exchange symmetries and branching rules for which such patent intricacy may eventually require groupoid techniques. An alternative approach to Lie algebra representations due to Jiménez-Montaño et al (1996) is to consider representations on hypercubes as based on Gray code structures. Already some known group structures show up here for various assortments of codon doublets and since sub-symmetries of these representations involve cubical methods, patterns of groupoid symmetries are also likely to arise. Thus we approach increasingly complex situations involving *groupoid representations* (see e.g. Bos, 2010) and *groupoid symmetry breaking*, techniques that can be computationally highly non-trivial, since even for relatively straightforward symmetries such as those appearing in certain ‘windmill patterns’, constraints do apply in order to facilitate programming capabilities (see e.g. Gent et al., 2010). Other questions may open up such as the possibility of breaking ‘mirror symmetry’ states in the genetic code caused by biochemical perturbations of chiral fields at the molecular level (Avetisov and Goldanskii, 1996).

### 3.4 Amino acid encoding–codon decoding and error load

In order for free energy and error load to fit into the picture, we follow part of the set up of the error-correction network analysis of Tlustý (2007, 2008) (cf Sella and Ardell, 2002). We take an amino acid  $\alpha$  to be encoded by a unique codon  $j$  represented in the encoder matrix  $[E_{\alpha j}]$  satisfying  $\sum_j E_{\alpha j} = 1$ , and similarly, the decoder matrix  $[D_{j\beta}]$  satisfying  $\sum_\beta D_{j\beta} = 1$ , means that each codon is translated into a unique amino acid  $\beta$ . Let  $\mathcal{N}_c$  denote the number of protein chains for  $c$  codons. Next we set

$$R_{ij} = P(\text{the probability that codon } i \text{ may be read correctly as or misread as } j), \quad (3.13)$$

and then let  $[R_{ij}]$  denote the reading matrix and  $C_{\alpha\beta}$  be the chemical distance between the original amino acid  $\alpha$  and the one that is read as  $\beta$ . On setting  $P_\alpha = P(\text{amino acid } \alpha \text{ is required})$ , the *error load*  $H_{\text{ED}}$  (the average distortion in an  $R(D)$  problem) of the map specified by *encoding/decoding* can be expressed in terms of paths  $P_{\alpha ij\beta}$ , specifically by

$$H_{\text{ED}} = \sum_{\alpha \rightarrow i \rightarrow j \rightarrow \beta} P_{\alpha ij\beta} C_{\alpha\beta} = \sum_{\alpha, i, j, \beta} P_\alpha E_{\alpha i} R_{ij} D_{j\beta} C_{\alpha\beta}. \quad (3.14)$$

This leads to a ‘take-over’ probability given by  $P_{\text{ED}} \sim \exp(-H_{\text{ED}}T^{-1})$  and to the *average error load*  $\langle H \rangle$  as follows. If we take  $S$  to denote entropy due to random drift, and  $T$  to be inversely proportional to average error size (the strength of the random drift relative to the selection force that pushes towards maximization), then this probability can be seen to minimize a functional analogous to the Helmholtz free energy  $F$  in terms of the average error load  $\langle H \rangle$  as in (3.5):

$$F = \langle H \rangle - TS = \sum_{\text{ED}} H_{\text{ED}} P_{\text{ED}} + T \sum_{\text{ED}} P_{\text{ED}} \ln P_{\text{ED}}, \quad (3.15)$$

which effectively averages out the difference between the genetic message relayed by a codon statement and that which is actually expressed by the genetic/epigenetic translation machinery itself.

## 4 Meaningful paths, robustness and error correction

### 4.1 Meaningful paths

We now specify our observations in a more general context. Suppose we consider a pattern of signalling input  $\mathbf{S}_i$  describing the state of the protein with initial codon stream  $\mathbf{S}_0$  to be mixed in an unspecified but systematic algorithmic manner with a pattern of an otherwise unspecified ongoing activity, including cellular, epigenetic and environmental signals  $\mathbf{W}_i$  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 in an iterative form according to which

$$\mathbf{S}_{i+1} = f([\mathbf{S}_i, \mathbf{W}_i]) = f(a_i), \quad (4.1)$$

for some unspecified function  $f$ . Comparing this with the situation in §4.2,  $\mathbf{S}$  would be a vector,  $\mathbf{W}$  a matrix, and  $f$  a product of their function at some time stage  $i$ . 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 as follows. Let

$$\begin{aligned} B_0 &\equiv b_0, \dots, b_k, \\ B_1 &\equiv b_{k+1}, \dots, b_m. \end{aligned} \quad (4.2)$$

Then:

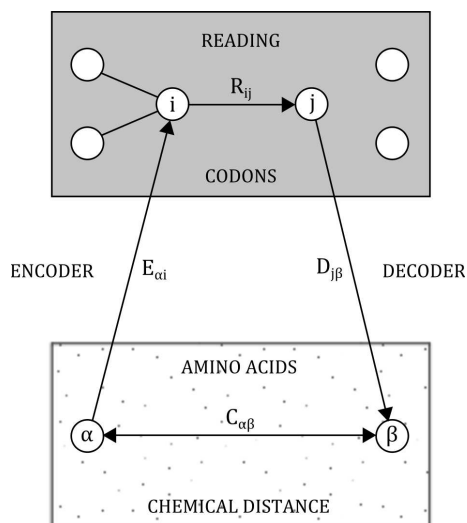


Figure 2: A schematic representation of the process in §3.4 as based on Fig. 2 of Tlustý (200a) where the error load gives the distortion of information described by  $H_{ED}$ .

(1) assume a graded response, supposing that if

$$h(x) \in B_0, \quad (4.3)$$

the pattern is not recognized, and

(2) if

$$h(x) \in B_1, \quad (4.4)$$

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

Expecting the coding signals to filtered appropriately (cf Ardell and Sella, 2002), we can further assume that  $B_0$  and  $B_1$  admit countable filtrations of the sort:

$$\begin{aligned} B_0 &= B_0^0 \subseteq B_0^1 \subseteq B_0^2 \subseteq \dots \\ B_1 &= B_1^0 \subseteq B_1^1 \subseteq B_1^2 \subseteq \dots \end{aligned} \quad (4.5)$$

where at level  $j$  we have set  $B_0^j \equiv b_0^j, \dots, b_k^j$ , and  $B_1^j \equiv b_{k+1}^j, \dots, b_m^j$ . Note that these oscillators may be influenced by ‘forcing’ when a signal is subjected to some impulse such that its frequency, and hence the response, adjusts accordingly with respect to an applied impulse. More familiar oscillating physical systems may react accordingly by exhibiting beats and resonance, for instance.

The principal objects of formal interest are paths  $x$  which, through information flow, trigger patterns of recognition-and-response. That is, given a fixed initial state  $a_0 = [\mathbf{S}_0, \mathbf{W}_0]$ , 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  $h(a_0, \dots, a_m) \in B_1$ . We can view  $B_1$  then as the set of final possible states  $\mathbf{S}_f \cup \{\mathbf{S}_{\text{path}}\}$  that includes both the final physical states and the set of all possible pathological conformations (see Wallace, 2010a, Figure 3).

For each positive integer  $n$ , let  $N(n)$  be the number of high probability grammatical/syntactical paths of length  $n$  which begin with some particular  $a_0$  and further leading to the condition  $h(x) \in B_1$ . These are paths of combined signals as above, that are structured to some language. For short, we 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  $h(x) \in B_1$ .

One critical assumption which permits an inference on the necessary conditions constrained by the asymptotic limit theorems of information theory, is that the entropy, as defined by the finite limit

$$H \equiv \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n}, \quad (4.6)$$

both exists and is independent of the path  $x$ . The rate distortion principle applies as follows (Wallace, 2005): *the restriction to meaningful sequences of symbols increases the rate at which information can be transmitted with arbitrary small error, and that the grammar/syntax of the path can be associated with a dual information source*. Besides the DMI and DMC properties, we may also assume a typical information source  $\mathbf{X}$  to be ‘adiabatic’, ‘piece-wise stationary’ and ‘ergodic’ (APSE), and that



the relevant systems engaging in a bio-cognitive process is describable as such. More specifically, the essence of ‘adiabatic’ is that, when the information source is parametrized according to some appropriate scheme, within continuous ‘pieces’ of that parametrization, alterations in parameter values occur slowly enough so that the information source  $\mathbf{X}$  remains as close to stationary and ergodic as necessary in order to implement the limit theorems of information theory. In this way ‘structure’ is subsumed within the sequential grammar and syntax of the dual information source, rather than within the sets of developmental paths as considered in Wallace and Wallace (2010).

In view of (4.6), the Shannon entropy of  $\mathbf{X}$  can be stated more specifically by (see e.g. Cover and Thomas, 1991):

$$H[\mathbf{X}] = \lim_{n \rightarrow \infty} \frac{\log[N(n)]}{n}. \quad (4.7)$$

With respect to e.g. the robustness criteria of §4.2, the time dependent information sources  $\mathbf{X}_i(t)$  are identified with the  $i$ -th component of the expressional pattern  $\mathbf{S}(t)$ , that is, we assign  $\mathbf{X}_i(t) \mapsto \mathbf{S}_i(t)$ , where as before  $\mathbf{S}_i(t) = f(a_{i-1})$ .

Recalling how the information source uncertainty was defined as in equation (4.6), an essential observation is a *fundamental homology* with the free energy density of a thermodynamical system such as that displayed in equation (3.6). Such a homology arises from Feynman’s observations (Feynman, 1996) reflecting in part on Bennett (1982) where this homology is effectively an identity, at least for very simple systems. From a more general perspective, Feynman (1996) postulates the information contained in a message as proportional to the amount of free energy density needed to erase it. This simply amounts to the fact that computing in any form takes work and the more complicated a coding or signalling process so measured by its source uncertainty, the greater its energy consumption. Putting it another way, the less information available to us concerning an event the higher is its entropy, and information retrieved is not without a cost in expenditure (of energy), where ‘cost’ is interpreted as the necessary number of bits needed to encode a message (the thermodynamic minimum of energy in terms of bits of information is  $k_B T \log_2 e$  erg/bit or  $k_B T$  erg/nat). So the efficiency in an information system essentially happens when there is the minimum amount of energy expended in retrieving information. Specifically, if  $F$  is taken to denote the free energy, then setting  $\Lambda$  equal to the minimum number of nats/sec, the efficiency of the system is given by  $\eta = k_B T F^{-1} \Lambda$  (see e.g. Berger, 1971).

## 4.2 Transcriptional regulators and robustness

There are certain evolutionary innovations resulting from an interplay of mutations and natural selections whereby, in a descriptive sense, a genotype corresponds to a regulatory network with a given topology and a phenotype to that of a steady state genetic pattern. This mechanism is constrained by certain conditions requiring processes to sustain a degree of robustness, meaning here a resilience towards environmental perturbations and thermodynamic effects, while at the same time admitting some ‘diversity’ in the process of messaging reception. Such a function of evolution and environment is to ensure that proteins can continue their catalyzing role in the presence of amino acid mutations, that the regulatory networks can continue to function in a noisy environment, and that embryos can develop normally in the presence of such perturbations. In any case, these regulatory networks, (protein) synthesis and the mutational operations can be seen as part and parcel with the question of folding (misfolding), while observing that error-minimization permits the appropriate codon allocation to amino acids through sequences of broken symmetries in terms of tRNA mutations (see Bashford et al., 1998, 2008).

Thinking back to the context of §4.1, we next turn to an analogous, but closely related sequence of  $N$  transcriptional regulators represented by their expressional patterns  $\mathbf{S}(t) = (\mathbf{S}_1(t), \mathbf{S}_2(t), \dots, \mathbf{S}_N(t))$ , in network form, at some time  $t$  that can influence expressions between themselves via cross-regulatory and auto-regulatory interactions as expressed by a matrix  $W = [w_{ij}]$  where  $w_{ij}$  represents a signaled regulatory influence  $w_{ij} : \text{gene } i \Rightarrow \text{gene } j$ , given the rules (1)  $w_{ij} > 0$  means activating, (2)  $w_{ij} < 0$  repressing, and (3)  $w_{ij} = 0$  absence.

In Ciliberti et al. (2007) such regulatory interactions describe the expressional state of the network  $\mathbf{S}(t)$  akin to a typical spin-glass model (Sourlas, 1989, see also Appendix §C), as specified by

$$\mathbf{S}_i(t + \tau) = \sigma \left[ \sum_{j=1}^N w_{ij} \mathbf{S}_j(t) \right], \quad (4.8)$$

where  $\tau$  is a constant and  $\sigma(\cdot)$  is a sigmoidal function  $\sigma : \mathbf{S}(t) \rightarrow (-1, 1)$ . For instance, with strong cooperation we may have  $\sigma = \text{sgn}$ , giving  $\mathbf{S}_i = \pm 1$ . Here  $\mathbf{S}(t)$  can be taken as an *incoming input*, mixed in a systematic way relative to  $W = [w_{ij}]$ , to create a path of combined signals  $x = (a_0, a_1, \dots, a_n, \dots)$  as to be seen in §4.1 homologous to the sequence  $\mathbf{S}(t + \Delta t)$ , with  $n = t(\Delta t)^{-1}$ , where on recalling expression (4.1) we set  $\mathbf{S}_{i+1} = f([\mathbf{S}_i, \mathbf{W}_i]) = f(a_i)$ . Accordingly, the structure becomes as much of a function of the sequential grammar and syntax of the dual information source as it is the cross-sectional intervals of the space of the  $W = [w_{ij}]$  (see Wallace, 2009) where, for instance, one would denote by  $\mathbf{S}(0)$  an initial state and by  $\mathbf{S}_\infty$

a stable equilibrium state, with a distance measure  $\mathcal{D}$  for graph topologies  $W, W'$  taken to be

$$\mathcal{D}(W, W') = \frac{1}{2M_+} \sum_{i,j} |\text{sgn}(w_{ij}) - \text{sgn}(w'_{ij})|, \quad (4.9)$$

where  $M_+$  denotes the maximum number of regulatory interactions.

In essence this construction reveals that genotype space, for instance, can be traversed in small increments without changing the phenotype which has evolutionary significance for genetic patterns: randomly selected pairs of networks of the same phenotype may have very different structure and may be subject to varying selective pressures. Thus one may imagine that a large overall ‘diameter’ of the network may be a critical feature for diversity of phenotype, and because some lengthy travel across the graph may be necessary to find all new phenotypes Ciliberti (2007), a distance measure of two phenotypes  $\mathbf{S}_\infty, \mathbf{S}'_\infty$  is given by the Hamming distance  $d_H$  in the form

$$d_H = d_H(\mathbf{S}_\infty(j), \mathbf{S}'_\infty(j)) = 1 - \sum_j \frac{\delta}{N} [\mathbf{S}_\infty(j) - \mathbf{S}'_\infty(j)], \quad 0 \leq d_H \leq 1, \quad (4.10)$$

where Kronecker  $\delta = 1$  should both arguments be equal and  $\delta = 0$  otherwise. Note that for such Hamming codes it is a basic fact that decoding all patterns of length  $\leq k$  is equivalent to  $(d_H)_{\min} \geq 2k + 1$  (see e.g. McElice, 2004; Yockey, 2005).

Related is how, in the statistical mechanics formulation, genetic algorithms based on spin glass models can reveal optimal selectivity as increasing with evolution. Prügel-Bennett and Shapiro (1994), show how selecting those solutions that are at a higher level of fitness, can be paired (through a crossover operation say) and then tested. This is performed iteratively through an algorithm up to the point where there is no further improvement in the examined population. Using spin glass states, Prügel-Bennett and Shapiro (1994) apply a chain as represented by vectors of the spins  $\sigma^{(\alpha)}$  (where  $\alpha = 1, \dots, P$ ) indexed by different members of the population; this spin vector is implemented in the genetic algorithm. Typically, new spins  $\tau_i^\alpha = \sigma_i^\alpha \sigma_{i+1}^\alpha$  are created. Selectivity on the basis of mutation and crossover follows from the energy levels of the Ising spin glass (which is described here in Appendix C).

## 5 Rate distortion coevolutionary dynamics

### 5.1 The basic equations

Understanding the time dynamics of cognitive systems away from phase transition critical points thus requires a phenomenology similar to the Onsager relations. 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 view of the analogy with nonequilibrium thermodynamics, the gradients in the  $K_j$  of the *disorder*, defined as

$$S \equiv H(\mathbf{K}) - \sum_{j=1}^m K_j \partial H / \partial K_j, \quad (5.1)$$

are of central interest. Note that equation (5.1) is analogous to the definition of entropy in terms of the free energy density of a physical system, as suggested by the homology between the latter and the information source uncertainty. Pursuing the homology further, the generalized Onsager relations defining temporal dynamics become

$$dK_j/dt = \sum_i L_{ji} \partial S / \partial K_i, \quad (5.2)$$

where the kinetic coefficients  $L_{ji}$  are, in first order, constants interpreted as reflecting the nature of the underlying cognitive phenomena (without requirement of the symmetry condition  $L_{ij} = L_{ji}$ ). The partial derivatives  $\partial S / \partial K$  are analogous to thermodynamic forces in a chemical system, and may be subject to override by external physiological driving mechanisms as shown in Wallace (2005), Wallace and Fullilove (2008) along with further extensions of these dynamical procedures.

Induced by the fundamental homology between the Shannon entropy and free energy density, the rate distortion  $R(D)$  follows a homologous path relation to the latter, thus suggesting that the dynamics of any bio-cognitive module interacting in characteristic real-time  $\tau$  will be constrained by the system as described in terms of  $R(D)$ . This can be seen more generally (Wallace and Wallace, 2008, 2009) by producing a vector-valued function  $R(\mathbf{Q})$  where in the vector  $\mathbf{Q} = (Q_1, \dots, Q_k)$  the first component is defined to be the average distortion, and then (cf (5.1)), we have

$$S_R \equiv R(\mathbf{Q}) - \sum_{i=1}^m Q_i \partial R / \partial Q_i, \quad (5.3)$$

which leads to the deterministic and stochastic systems of equations analogous to the Onsager relations of nonequilibrium thermodynamics

$$dQ_j/dt = \sum_i L_{ji} \partial S_R / \partial Q_i, \quad (5.4)$$

together with

$$dQ_t^j = L^j(Q_1, \dots, Q_k, t) dt + \sum_i \sigma^{ji}(Q_1, \dots, Q_k, t) dB_t^i, \quad (5.5)$$

where the  $dB_t^i$  represents often highly structured stochastic noise whose properties may be described in terms of Brownian motion and quadratic variation (see e.g. Protter 1995).

## 5.2 Phenomenological Onsager relations

Here we turn to different developmental subprocesses of gene expression characterized by information sources  $H_m$  interacting via chemical or other types of signals, and assume that different processes become each other's principal environments. This is a working hypothesis within a broad coevolutionary context that underscores the cognitive element. Let

$$H_m = H_m(K_1, \dots, K_s, \dots, H_j, \dots), \quad (5.6)$$

where the  $K_s$  represent other relevant parameters, and  $j \neq m$ . We regard the dynamics of this system as driven by a recursive network of stochastic differential equations. Letting the  $K_j$  and  $H_m$  all be represented as parameters  $Q_j$  (with the caveat that  $H_m$  does not depend on itself), we follow the generalized Onsager formulation of Wallace and Wallace (2009), in terms of the equation

$$S^m = H_m - \sum_i Q_i \partial H_m / \partial Q_i, \quad (5.7)$$

to obtain a recursive system of *phenomenological Onsager relations*, in terms of a system of stochastic differential equations

$$dQ_t^j = \sum_i [L_{ji}(t, \dots, \partial S^m / \partial Q^i, \dots) dt + \sigma_{ji}(t, \dots, \partial S^m / \partial Q^i, \dots) dB_t^i], \quad (5.8)$$

in which, for ease of notation, both the terms  $H_j$  and the external  $K_j$ 's are expressed by the same symbol  $Q_j$ . As  $m$  ranges over the  $H_m$  we could allow different kinds of 'noise'  $dB_t^i$ , having particular forms of quadratic variation which may represent a projection of environmental factors within the scope of a rate distortion manifold (Glazebrook and Wallace, 2009b). The noise factor is significant in view of the findings of Austin et al. (2006) who have observed that perturbations of the network parameters inducing stochastic fluctuations in the molecular patterns may in turn influence regulatory mechanisms, and in similar way to how the presence of stochastic resonance may amplify certain signals, noise-spectral measurements may then uncover further mechanisms which could be potentially beneficial to the code's evolution.

We remark that equation (5.8) can be generalized somewhat (Wallace and Wallace, 2009) with respect to crosstalk, its distortion, the inherent time constants of the various bio-cognitive modules, and in particular, the overall available free energy density. As shown in Glazebrook and Wallace (2010), analysis of the rate distortion dynamics on a case-by-case basis, motivates integration to a multidimensional Itô process as given by

$$Q_t^\alpha = Q_0^\alpha + \sum_{\beta=\{ij\}} \left[ \int_0^t L_\beta(s, \dots, \partial S_R^\beta / \partial Q^\alpha, \dots) ds + \int_0^t \sigma_\beta(s, \dots, \partial S_R^\beta / \partial Q^\alpha, \dots) dB_s^\beta \right], \quad (5.9)$$

and this in turn leads to a stochastic flow on a suitable topological manifold which in this present context could serve as a more general model for codon space. In fact, such a flow property had already been observed in Tlustý (2007), namely, that the standard genetic code and its variants evolve as a flow within the codon space. However, given that "freezing" of some sort is likely to re-occur in the quest for optimal error-correction, we expect such a flow to be stalled at certain time intervals, thus creating singularities in the flow in a dynamical systems sense (an analytic technicality to be finessed here).

## 5.3 Metric on a space of languages

Let us note that equations (5.1) and (5.2) can be derived in a simple parameter-free covariant manner which relies on the underlying topology of the information source space that is implicit to the processes which we envisage. Different bio-cognitive phenomena have, according to our development, dual information sources, and we 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 bio-cognitive process, there is an open set  $U$  of closely similar languages  $\hat{A}$ , such that  $A$  and  $\hat{A}$  are subsets of  $U$ .

Since the information sources dual to the processes are similar, for all pairs of languages  $A, \hat{A}$  in  $U$  within a given embedding alphabet, we define a metric on the latter by

$$\mathcal{M}(A, \hat{A}) = \left| \lim \frac{\int_{A, \hat{A}} d(Ax, \hat{A}x)}{\int_{A, A} d(Ax, A\hat{x})} - 1 \right|, \quad (5.10)$$

with respect to a distortion measure  $d(Ax, \hat{A}x)$ , and apply standard integration arguments over the high probability paths, where the usual metric properties apply as in e.g. Burago (2001). In the context of Ardell and Sella (2002), we may see such a metric as derived from an informational driven physico-chemical distance function with respect to the analogous  $A$  and  $\hat{A}$  coding. Also, since  $H$  and  $\mathcal{M}$  are both scalars, a ‘covariant’ derivative can be defined directly as

$$dH/d\mathcal{M} = \lim_{\hat{A} \rightarrow A} \frac{H(A) - H(\hat{A})}{\mathcal{M}(A, \hat{A})}, \quad (5.11)$$

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

A relatively straightforward case is the following. Suppose the system is set in some reference configuration  $A_0$ . To obtain the unperturbed dynamics of that state, impose a Legendre transform using this derivative, defining another scalar

$$S \equiv H - \mathcal{M}dH/d\mathcal{M}. \quad (5.12)$$

Then simplest possible Onsager relation – here seen as an empirical, fitted, equation like a regression model, becomes

$$d\mathcal{M}/dt = LdS/d\mathcal{M}, \quad (5.13)$$

where  $t$  is the time and  $dS/d\mathcal{M}$  represents an analog to the thermodynamic force in a chemical system (cf §6.4 of Berger, 1971).

## 5.4 Mutations: mutual entropy between sequence-structure

As analogous to the expressional patterns of §4.2, the previous techniques are applied to the following case of mutations which are themselves functions of evolution, and together with selection and translational error, can influence the distribution of codons to the extent that the latter favor patterns of error-correction that drift to some optimal level and can ameliorate mutation effects (Ardell and Sella, 2002; Sella and Ardell, 2002). For instance, let us consider a series of amino acid sequences

$$\{\dots, \text{Seq}_{t-1}, \text{Seq}_t, \text{Seq}_{t+1}, \dots\} = \{\text{Seq}_t\}_{t \in \mathbb{Z}}, \quad (5.14)$$

where each  $\text{Seq}_t$  applies to one protein chain, ordered by a discrete temporal order  $t \in \mathbb{Z}$  of corresponding tertiary structures

$$\{\dots, \text{Str}_{t-1}, \text{Str}_t, \text{Str}_{t+1}, \dots\} = \{\text{Str}_t\}_{t \in \mathbb{Z}}. \quad (5.15)$$

Such a chain can be represented as a noisy digital communication channel and with an output probability of at least  $\sim 30\%$  with a Shannon limit at  $10^{-2}$  bits/amino acid (see Lisewski, 2008). Ardell and Sella (2002) claim that codes evolving with messages that mutate under such a process, tend to freeze with redundancy. This can be reduced to analyzing three different possibilities: the coevolution of genetic codes with:

- (1) transitional-biased message mutation and no translation misreading;
- (2) translational misreading and no transition bias in mutation;
- (3) transition-biased message mutation and translational misreading.

An example in Lisewski (2008) considers concatenated primary sequences  $\{\text{Seq}_t\}_{t \in \mathbb{Z}}$  resulting in a stream of letters from the amino acid alphabet  $A$  with (alphabetical) size  $|A| = 20$ . The encoder is a map that uses a block code of fixed length  $n$ , say, to encode the source through the code book; in other words, a map for every sequence

$$\text{Seq}_t \longrightarrow (\text{single code word})X^n(\text{Seq}_t), \quad (5.16)$$

represented by an  $n$ -vector  $(X_1, \dots, X_n)$  of integers. The code word in turn belongs to the book of 20 possible structure symbols  $A^* = \{a_1^*, \dots, a_{20}^*\}$ , the finite set of all code words corresponding to the 20 amino acid symbols  $\{A, G, \dots\}$ , where  $a_j^* \in A^*$  are contact vectors determining the amino acid sequence. The message input term  $X^n(\text{Seq}_t)$  from (5.16) is relayed over a

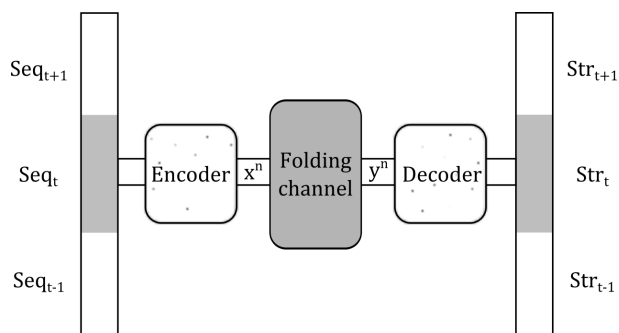


Figure 3: Simplified diagrammatic representation of the sequence-structure information process of §5.4 as based on Fig. 1 of Lisewski (2008).

noisy channel which then outputs an  $n$ -vector  $\Upsilon^n(\text{Str}_t) = (Y_1, \dots, Y_n)$  representing the folded protein chain  $\text{Str}_t$  following which a single use of the channels is the transmission of a single amino acid sequence subject to the *channel capacity*

$$C = \max_{p(A)} I(A, A^*). \quad (5.17)$$

In view of §5.3, we modify the role of  $\hat{A}$  via the assignment  $\hat{A} \mapsto A^*$ , and for times stages  $t, t'$ , take as above, the metric  $\mathcal{M}(\text{Str}_t, \text{Str}_{t'})$ . At each side of the communication channel we have for the symbol sequences  $|S_A| = 7702314$  amino acid symbols and  $|S_{A^*}| = 31609$  corresponding structural symbols (Lisewski, 2008).

As for the code rate, we have  $R(D) = H(A)/n$ , where  $H(A)$  is interpreted as the Shannon entropy of the amino acid sequence, where  $n$  is the code block length implemented by the encoder. Assuming the code rate  $R(D)$  and channel capacity are known, then in accordance with the Rate Distortion theorem, we have  $R(D) < C$ , leading to, for every block size,  $n > n_{\min} = H(A)/C$ , and the codes exist, and no such code when  $R(D) \geq C$ . The Shannon entropy  $H(A) = 3.90$  bits for the amino acid alphabet  $A$ , and  $H(A^*) = 3.76$  bits for the structural code words in  $A^*$  (Lisewski, 2008). Further, the mutual entropy between structure and sequence following Adami (2004) is given by

$$I(\text{Seq}_t : \text{Str}_t) = H(\text{Seq}_t) - H(\text{Seq}_t | \text{Str}_t), \quad (5.18)$$

and should the environment directly influence the structure, then we would have

$$H(\text{Str}_t | \text{Seq}_t) \simeq H(\text{Seq}_t | \text{Env}_t). \quad (5.19)$$

When taking  $H(\text{Str}_t | \text{Seq}_t) = 0$ , we can re-formulate (5.18) as

$$\begin{aligned} I(\text{Seq}_t : \text{Env}_t) &\simeq I(\text{Seq}_t : \text{Str}_t) \\ &= H(\text{Str}_t) - H(\text{Str}_t | \text{Seq}_t) \\ &= H(\text{Str}_t), \end{aligned} \quad (5.20)$$

which in view of the mutual entropy between sequence and structure, expresses how the thermodynamical entropy of possible protein structures can be constrained by information about the environment as it is coded by the sequence. For instance, excessive noise and random inputs of symbols in  $S_{A^*}$  would most probably corrupt a corresponding code in  $A^*$ , and once again the Shannon estimate serves as a threshold should errors exceed a critical bound. On the empirical side the Protein Data Bank (PDB) provides sequence-structure data giving  $H(A) = 3.90$  bits, with block length  $n = 400$ , with transmission rate  $R(D) = 0.010$  bits per amino acid symbol followed, with channel capacity estimated at  $C = 0.016$  bits (per amino acid symbol). When restricted to  $\mathcal{N}_{25} = 2372$  protein chains with mutual sequence identity of  $< 0.25$ , the estimated  $C_{(25)} = 0.016$  bits, was attained (see Lisewski, 2008, Figure 4.).

## 6 The Topological Hypothesis, phase transitions

### 6.1 The codon space as a graph

The carrier for the dynamics surveyed here is modeled on a rate distortion manifold which has wide-scale overlap with those codon spaces structured in a way that evolution can be influenced by mapping out those regions which can accommodate

Code	# Codons	Max. # AA's
4-base singlets	4	4
3-base doublets	9	7
4-base doublets	16	11
16 codons	32	16
48 codons	48	20
4-base triplets	64	25

Table 1: Part of Table 1 in Tlustý (2007a) showing the topological limit to the number of amino acids for different codes.

load minimization and diversification so that site type, coding fitness, targets, etc. can be correlated (as in Ardell and Sella, 2002). One expects the rate distortion manifold to have, in an analytic sense, some degree of differentiability, though here we will finesse this issue and choose to consider underlying combinatorial structure. Specifically, we let  $\Gamma = (V, E)$  denote a graph with  $V$  denoting a finite vertex set,  $E$  an edge set with an oriented edge  $e = (u, v)$  (accordingly,  $e^{-1} = (v, u)$ ) such that  $u = i(e)$  is the initial vertex and  $v = t(e)$  is the terminal vertex, and let  $F$  be the number of enclosed faces. Tlustý (2007a, 2008) formulates the code that emerges at the phase transition appears a mode  $e_{\alpha i}$  that minimizes the free energy  $F$ . The codon space can be described as such a graph  $\Gamma$  whose vertices are the codons and two codons  $i, j$  are linked by an edge if (see §3.4) there exists an associated  $R_{ij} (\neq 0)$  in the reading matrix, under the following conditions/observations:

- (1) The vertex set  $V$  consists of codons whereby two codons are linked by an edge in the likelihood they may be confused by misreading.
- (2) Two codons are most likely to be confused if all their letters except for one agree and then they are connected by an edge. The resulting graph  $\Gamma$  is natural for considering the impact of translation errors on mutations because such errors almost always involve a single letter difference, that is, a movement along an edge of the graph to a neighboring vertex.
- (3) The native state of the protein has the lowest available free energy induced by the interaction of the amino acid sequence with the embedding environment.
- (4) Recall that there is an embedding  $\Gamma \rightarrow S$  into a surface  $S$  and the topology of  $\Gamma$  is characterized by its genus  $\gamma(S)$  which is the minimal number of holes required for  $\Gamma$  to be embedded in  $S$  such that no two edges cross. For the underlying network we have the well-known combinatorial formula  $\gamma = 1 - \frac{1}{2}(V - E - F)$ .

Thus the greater the number of connected components in the graph, the higher the genus becomes for a minimal embedding. In Tlustý (2007a) the interconnected 64-codon graph can be embedded in a surface with genus  $\gamma(S) = 41$ . If only 48 effective codons are considered, then the genus is reduced to  $\gamma(S) = 25$ .

In view of this criteria, the claim is that the evolution of the code is determined by the underlying topology of its graph and in a transitional phase, it is only those modes with the least error-bound that can emerge and are subjected to alteration by the topology. From the perspective of Pettini (2007), a free energy argument serves as a Morse function whose critical points characterize just such a topology. Tlustý (2007) in a more specific way considers the topology of the code as imposing an upper limit to the number of low modes – critical points – of the corresponding free energy-analog functional, and this is also the number of amino acids. The low modes define a partition of the codon surface into domains, and in each domain a single amino acid is encoded. The partition optimizes the average distortion by minimizing the boundaries between the domains as well as the dissimilarity between neighboring amino acids. This bound on the number of low nodes (and thus as claimed, the number of amino acids) arises as an application of the well-known *chromatic number* as given by Heawood's formula (see Ringel and Youngs, 1968):

$$\text{chr}(\gamma(S)) = \text{int}\left[\frac{1}{2}(7 + \sqrt{1 + 48\gamma(S)})\right], \quad (6.1)$$

where  $\text{chr}(\gamma(S))$  is the number of color domains of a surface  $S$  with genus  $\gamma(S)$ , and  $\text{int}[x]$  denotes the integer value of  $x$ . Recall also that the Euler characteristic  $\chi(S) = 2 - 2\gamma(S)$ . In particular, in Tlustý (2007a, 2008) it is the genus that represents the number of holes in the protein folding error network associated with the code and the chromatic number  $\text{chr}(\gamma(S))$  is a measure of the number of protein symmetries.

More generally, for a topological manifold  $M$  having a Morse function  $F$ ,  $\chi(M)$  can be expressed as the alternating sum of the function's Morse indices  $\mu_i (i = 0, 1, \dots, m)$  of  $F$  on  $M$ , defined as the number of critical points ( $dF(x_c) = 0$ ) of index  $i$ , that is, the number of negative eigenvalues of the matrix  $H_{i,j} = \partial^2 F / \partial x_i \partial x_j$ . Then by the Poincaré-Hopf theorem,

$$\chi(M) = \sum_{i=0}^m (-1)^i \mu_i, \quad (6.2)$$

genus $\gamma(S)$	$\text{chr}(\gamma(S))$
0	4
1	7
2	8
3	9
5	10
6,7	11
8,9	12

Table 2: From Wallace (2010c) where the genus  $\gamma(S)$  equals the number of network holes, and using the formula (6.1) the chromatic number  $\text{chr}(\gamma(S))$  gives the number of protein symmetries in each case

which holds true for any Morse function on the manifold  $M$  (see e.g. Matsumoto, 2001, and Appendix §B.2 here).

**Remark 6.1.** Applying a spontaneous symmetry lifting argument to  $F_R$  generates topological transitions in codon graph structure as the ‘temperature’  $R(D)$  increases, i.e., as the average distortion  $D$  declines, via the inherent convexity of the Rate Distortion Function. That is, as the channel capacity connecting codon machines with amino acid machines increases, more complex coding schemes become possible. In this respect, we recall that for the surface  $S$ , the Euler characteristic  $\chi(S) = 2 - 2\gamma(S)$  as in (B.4) can be expressed in terms of the cohomology structure of  $S$  (e.g., Lee, 2000, Theorem 13.38) where by the Poincaré Duality Theorem, the homology groups of a manifold are related to the cohomology groups in the complementary dimension (e.g. Bredon, 1993, p. 348) and thus points to the ‘fundamental homology’ described earlier. One can then envisage the (co)homology groupoid to be taken as the disjoint union of the (co)homology groups of the embedding manifold.

## 6.2 Spectrum of the graph Laplacian

Next we consider the Laplacian  $\Delta$  of  $\Gamma$ . If a pair of vertices  $(i, j) \in E$  are adjacent, then in terms of e.g. the reading matrix  $[R_{ij}]$  (with  $R_{ij} > 0$ ), we have

$$\Delta_{ij} = \Delta_{ji} = -R_{ij} < 0, \quad (6.3)$$

otherwise  $\Delta_{ij} = 0$ , and  $\Delta_{ii} = -\sum_{i \neq j} \Delta_{ij}$  (see Appendix §B.3). For instance, if  $\Gamma$  is taken to be the error graph of §6.1, then  $\Delta$  is the operator that measures the effect of errors and so regulates any phase transition.

Corresponding to the  $n$ -th eigenvalue  $\lambda_n$ , the eigenfunction  $u_n$  admits at most  $n$  weak sign graphs; in particular, for  $n = 2$ , the eigenfunction  $u_2$  divides  $\Gamma$  into precisely two weak sign graphs (see §B.3). Thus it is of interest to determine the dimension of the corresponding eigenspace and multiplicity  $m$  of  $\lambda_2$ . The quantity  $m$  is a measure of the first energy excitation being the primal mode for types of continuous (or second order) phase transitions. The chromatic number  $\text{chr}(\gamma(S))$  of (6.1) identifies the maximal number of first excited modes of the  $\Delta$ . Letting  $\bar{m}(S)$  denote the supremum of  $m$  over all possible  $\Delta$  on  $S$ , there is the estimate of Colin de Verdière stating that  $\bar{m}(S) \geq \text{chr}(\gamma(S)) - 1$  (see e.g. Tlustý, 2007b). In the case of functions, the graph  $\Gamma$  is a reliable ‘spectral’ model for  $S$  in the sense that from Theorem 5.7 of Dodziuk (1976), the eigenvalues of all orders of  $\Delta$  on  $\Gamma$  converge to those of the continuous Laplacian on functions as defined on  $S$  (see B.3).

## 6.3 Phase transitions and holonomy

Given the graph  $\Gamma = (V, E)$ , the *star of a vertex*  $\text{st}(v)$  is the set of edges emanating from  $v$ , that is

$$\text{st}(v) = \{e : i(e) = v\}. \quad (6.4)$$

The various components of the graph may be thought of a comprising a cell network in which the coupling and equivalence of cells leads to a natural groupoid structure having a system of specific equivalence classes  $[v]_V$  and  $[e]_E$ , for vertices and edges, respectively (see Appendix §A.1). With the inclusion of this extra structure we then append  $\Gamma$  to  $\Gamma = (V, E, \sim_v, \sim_e)$ . Here the vertices (nodes) of the network are representative of certain cells where the synchrony of the system depends on groupoid symmetries which in a sense is broken by an impinging rapid crosstalk internal to the system while the latter attempts to manage a slower external crosstalk.

Next, we implement some general procedures based upon the idea of a *connection*  $\nabla$  on  $\Gamma$ , relative to the stars ( $\text{st}$ ) of vertices which following Bolker et al. (2006), is explained with some details in Appendix §B.1 as the combinatorial analog of ‘covariant differentiation’, a principle familiar to students of calculus. We take vertices  $(e_1, e_2, \dots, e_{k+1})$  interpreted as  $k + 1$

information sources  $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{k+1})$  in accordance with the APSE condition of §4.1, where the  $\mathbf{X}_i$  act with the set of tuning parameters. A connection  $\nabla$  is considered as an operation

$$\nabla(\mathbf{X}_i, \mathbf{X}_j) : \text{st}(\mathbf{X}_i) \longrightarrow \text{st}(\mathbf{X}_j), \quad (6.5)$$

for  $1 \leq i, j \leq k+1$ , satisfying certain properties (see §B.1). With respect to the metric  $\mathcal{M} = \mathcal{M}(\mathbf{X}_i, \mathbf{X}_j)$  applied to these information sources, the above connection in (6.5) implements on the underlying network, the covariant differentiation along the path  $\mathbf{X}_i \longrightarrow \mathbf{X}_j$ , just as in (5.11):

$$dH/d\mathcal{M} = \lim_{\mathbf{X}_j \longrightarrow \mathbf{X}_i} \frac{H(\mathbf{X}_j) - H(\mathbf{X}_i)}{\mathcal{M}(\mathbf{X}_i, \mathbf{X}_j)}. \quad (6.6)$$

Corresponding to each  $\mathbf{X}_i$ , a maximized channel capacity  $\mathbf{C}_i$  is assigned, in accordance with the Shannon estimate  $H(\mathbf{X}_i) \leq \mathbf{C}_i$ , for  $1 \leq i \leq k+1$ , thus respecting the Rate Distortion theorem along paths  $\mathbf{X}_j \longrightarrow \mathbf{X}_i$ . If necessary, we can view  $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{k+1})$  as comprising a closed geodesic, and as explained in §B.1, the set of these in a given graph will thus specify  $\nabla$ . Once we have a handle on  $\nabla$  it is then possible to apply to  $\Gamma$  certain operations analogous to the differential-geometric setting in order to explore the structural geometry of the various graphs which have been described so far (cf Glazebrook and Wallace, 2009a).

This technique of the network geometry can be applied to the entropy rates occurring in the various cases we have considered so far. For the sequence-structure-environment in the noisy communication channels with the data of §5.4, we assign  $\text{Str}_t$  (at time  $t$ ) to a corresponding sensory input  $\mathbf{S}_t$ , further combined with environmental signals  $W_t$  and combined signals  $a_t$  just as in (4.1):

$$\begin{cases} \text{Str}_{t+1} & = f([\text{Str}_t, W_t]) = f(a_t) \\ I(\text{Seq}_t : \text{Env}_t) & = H(\text{Str}_t) \end{cases} \quad (6.7)$$

(here we have made replacements  $i \mapsto t$  and  $j \mapsto t'$ ), where we make a straightforward assignment from the vertex information source, at time  $t$ .

$$\mathbf{X}_t \mapsto \text{Str}_t, \quad (6.8)$$

(and likewise for  $\text{Seq}_t$ ). Using the principle of (6.5) applied to the mutual information  $I(\text{Seq}_t : \text{Env}_t) = H(\text{Str}_t)$  in (5.20), leads to considering the covariant derivative

$$dH/d\mathcal{M} = \lim_{\text{Str}_t \longrightarrow \text{Str}_{t'}} \frac{H(\text{Str}_t) - H(\text{Str}_{t'})}{\mathcal{M}(\text{Str}_t, \text{Str}_{t'})}, \quad (6.9)$$

as representing the graph connection

$$\nabla(\text{Str}_t, \text{Str}_{t'}) : \text{st}(\text{Str}_t) \longrightarrow \text{st}(\text{Str}_{t'}), \quad (6.10)$$

where again at each time stage  $t$ , the Shannon estimate  $H(\text{Str}_t) \leq C_t$  is observed. Likewise, the error load  $H_{\text{ED}}$  of §3.4 expressed in terms of paths  $P_{\alpha ij\beta}$  in (3.14) and their concatenation, now become the meaningful paths of §4.1. In this present graph formalism these paths are considered as determined by edges  $e_\nu \in E$  where each  $\nu = \nu(\alpha ij\beta)$  is a multi-index of the path subscripts.

A property of the connection  $\nabla$  in (6.5) is its *holonomy* which can be best described by considering how in the traditional differential-geometric sense, a smooth connection implements the parallel translation of vectors around closed paths, and the induced representation of the space of the latter into a group of global symmetries is essentially the holonomy (of the connection). The classic example is the *Poincaré first-return map* of a dynamical system that incorporates typical phase transitions. In the combinatorial setting of Bolker et al. (2006), the holonomy of  $\nabla$  can be described formally in terms of permuting the ‘stars’ of vertices towards a *spatiotemporal reorientation*, as follows. Let  $\mathcal{C} = \{e_1, \dots, e_n\}$  be any cycle in the graph  $\Gamma$ , for which the terminal and initial vertices satisfy  $t(e_\alpha) = i(e_{\alpha+1})$  modulo  $n$ . Then the connection around  $\mathcal{C}$  leads to a permutation

$$\nabla_{\mathcal{C}} = \nabla_{e_n} \circ \dots \circ \nabla_{e_1} \circ \nabla_{e_0}, \quad (6.11)$$

of the star set  $\text{st}(u)$ . The *holonomy group*  $\text{Hol}(\Gamma, \nabla)_u$  at a vertex  $u$  of  $\Gamma$ , is the subgroup of the *permutation group of st}(u) generated by the permutations  $\nabla_{\mathcal{C}}$  over such all cycles  $\mathcal{C}$  that pass through the vertex  $u$ . A phase transition may then be represented by a permutation through vertices in  $\Gamma$  and such a ‘geometric phase’ accounts for how the various bio-cognitive modules shift gear and create a reorientation of the system.*

Now let us return to equivalence classes and the role of groupoids. This implements the above permutation groups of  $\text{st}(u)$ . A *holonomy groupoid* is obtained via the disjoint union

$$\text{Hol}(\Gamma, \nabla) = \bigvee_{u \in \Gamma} \text{Hol}(\Gamma, \nabla)_u, \quad (6.12)$$



which pieces together the local operations, and at the same time produces an equivalence class representation of the phase transition and its internal amplitudes. We summarize this as follows: *the holonomy groupoid represents a globalization of the local dynamic iterates by providing what is essentially a representation of the graph's path components onto some prevailing group of symmetries*. In the presence of symmetry breaking, it would be reasonable to consider the groups  $\text{Hol}(\Gamma, \nabla)_u$  as commensurable to some degree with, for instance, the corresponding Lie groups featuring in the  $\mathfrak{sp}(6)$  chain in (3.12), or that of the  $\mathfrak{sl}(6, 1)$  chain enumerated in Bashford et al. (1998, 2008).

## 7 Discussion and conclusions

The code's development passed through "accidental phases" created by probabilistic events that could be both regulated and manipulated by an evolving error-correction mechanism. Here we have viewed the latter within the framework of Shannon entropy in the context of the fundamental homology with the free energy density of a thermodynamical system. A common thread to this and other works suggests that increased selection forces may have been significantly enhanced by rate distortion dynamics in regard to the critical behavior of the free energy Morse function and varying topology, a function of which would have induced an order of redundancy so mandated by coevolution. Thermodynamic parameter changes in turn induced spontaneous symmetry breaking which we have shown can be captured by several techniques of representation theory. One can then invert Landau's arguments and apply them to the (co)homology groupoid in terms of the rising 'temperature'  $R(D)$ , to obtain a punctuated shift to increasingly complex genetic codes with increasing channel capacity. Our development here realizes mappings **codon space**  $\rightarrow$  **amino acid space** quite explicitly in the context of rate distortion manifolds.

On the other hand, the case of protein folding from an amino acid string is not an entirely random process, but may be the result of an evolved structured statement by an information source's uncertainty and the occurrence of mutations which may not have been all random but were subject to environmental forces. In view of the redundancy issue, the corresponding evolutionary processes may be capable of extending the code's expression from 20 to 25 amino acids with the possibility of there being many other protein folding codes (Tlusty, 2007a; cf Hornos and Hornos, 1994; Bashford et al., 1998, 2008). Having said this, we add that there still remain open questions concerning the role of  $R(D)$ , since this in turn drives punctuated changes in the genetic code and further exploration will be necessary. But what seems to follow from the collective processes we have described in explaining "the frozen accident", is that certain *adaptation effects* are in play (just as one finds in various neurocognitive and bio-sociological phenomena), and in this regard we quote from Ardell and Sella (2002):

... Our work has been motivated by the belief that the patterns of the standard genetic code may be explicable as adaptations of a system of information processing. If this turns out to be plausible and correct, we may say that adaptations have reduced the deleterious consequences of genetic and physiological error at a very fundamental level of biological organization ...

So the "frozen accident" by any reasonable account appears to have arisen as an evolutionary 'adaptation' against a temporary unreadiness (or an enforced over-robustness) to assimilate a barrage of highly intricate genetic messaging during which time error-correction patterns strived to evolve and crystallize accordingly in order to withstand ongoing selective pressures.

We also note that holonomy and symmetry breaking are mathematical concepts that arise essentially from the iterates of local-to-global processing and one such product of this is indeed the holonomy groupoid, a novel concept that has been introduced in this paper for the purpose of analyzing genetic networks. Further, the question of groupoid representations may uncover deeper conceptual issues in view of representation spaces that are spaces of operators ('fields of Hilbert or Banach spaces' as in e.g. Bos, 2010), a setting that may be compared the 'supersymmetric' model of Bashford et al. (1998, 2006), but one that is likely to be highly non-trivial and costly in a computational sense. Thus in view of the various methods we have brought to the forefront, we cannot fail to acknowledge the remarkable insight of E. Schrödinger who claimed that classical physics was insufficient for understanding fundamental life processes. In particular Schrödinger (1967) had envisaged the potential importance of information theory in evolutionary genetics, how living systems can be alterable under thermodynamic effects that are often the results of adverse biological contagion and that quantum mechanical effects might catalyze potential mutations, revealing the organization and evolutionary drive of the genetic code all the more extraordinary.

## A Appendix: Groupoids and their atlases

### A.1 Concept of a groupoid

Many bio-cognitive processes are naturally dynamical systems (see e.g. Glazebrook and Wallace, 2009a). One aim in these systems is to unify the internal and external symmetries, and to be able to reduce vast myriad-like network configurations into manageable schemes involving the corresponding equivalence classes analogous to those already mentioned in source

encoding/decoding, etc. in §3.1 (see also §5.3 below). A precise way of doing this lies within the categorical concept known as a *groupoid* (see e.g. Brown; 2006, Connes, 1994; Weinstein, 1996). In essence a groupoid  $G$  consists of both a set of objects  $X$  and a set of morphisms, or ‘arrows’, each of which project to an object in  $X$ , and all such morphisms are invertible.

**Remark A.1.** The most familiar example of a groupoid, as known to students of algebra, is that of a ‘group’ where there is a single object (‘the identity’). Hence groupoids can be viewed as extensions of the ‘group’ concept to sets of *multiple identities* thus providing a wide scope of applications to the dynamics of neurocognitive and socio-bioinformatic systems (see e.g. Glazebrook and Wallace, 2009a, and references therein).

A groupoid can be depicted by

$$\alpha, \beta : G \begin{array}{c} \xrightarrow{\alpha} \\ \xrightarrow{\beta} \end{array} X \quad (\text{A.1})$$

where the groupoid morphisms  $(\alpha, \beta)$  onto objects, are called the *range* and *source maps*, respectively. Informally, the groupoid represents a feature of built in reciprocity between its algebraic structures, internalizing and externalizing the prevailing symmetries. The morphisms  $\alpha, \beta$  satisfy certain algebraic relations of associativity, existence of two-sided identities, etc. (for details, see e.g. Brown; 2006, Connes, 1994; Weinstein, 1996). A groupoid can here be understood in relationship to a linkage by a meaningful path of an information source dual to a cognitive process for which the underlying principle is that: *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 to the same base point, and the 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.*

**Example A.1.** Since we have already mentioned equivalence classes in the context of source encoding/decoding, it seems appropriate to see how an equivalence relation  $\mathcal{R}$  defined on (a set)  $X$  takes shape as a groupoid. Here we have the two projections  $\alpha, \beta : \mathcal{R} \rightarrow X$ , and a product  $(x, y)(y, z) = (x, z)$  whenever  $(x, y), (y, z) \in \mathcal{R}$  together with an identity, namely  $(x, x)$ , for each  $x \in X$ . Moreover, the essential equivalence relations (classes) derived from a systems space (network) arise from the orbit equivalence relation of some groupoid  $G$  acting on that space (see e.g. Weinstein, 1996). In the context of connected (sub)networks/graphs with path concatenation, that are representable in terms of equivalence classes, natural groupoid structures arise in accordance with equivalence classes of relations  $\mathcal{R}(xy)$ , as above, that is simply interpreted as having an edge linking node  $x$  to node  $y$ . Conversely, a groupoid (of equivalence relations) admits an underlying graph structure via its implicit scheme of objects and morphisms between objects (for details, see e.g. Brown, 2006; Golubitsky and Stewart, 2006). Thus we have the two-way associations whereby ‘objects’ can be identified with ‘nodes’, and ‘morphisms’ identified with ‘edges’ in groupoids (of equivalence relations) and networks, respectively:

$$\begin{array}{ccc} \text{Network} & \xrightleftharpoons{\text{equivalence relation}} & \text{Groupoid} \\ \text{Network} & \xleftarrow{\text{underlying graph}} & \text{Groupoid} \end{array}$$

## B Appendix: Some geometry of the network architecture: geodesics and phase transitions

### B.1 Connections on graphs and geodesics

Firstly, for graph-theoretic models there are certain combinatorial notions which can be used to replicate a ‘differential’ structure as realized on a standard differentiable manifold (such as a sphere or a torus). Let  $\Gamma = (V, E)$  be a graph with  $V$  denoting a finite vertex set,  $E$  an edge set with an oriented edge  $e = (u, v)$  (accordingly,  $e^{-1} = (v, u)$ ) such that  $u = i(e)$  is the initial vertex and  $v = t(e)$  is the terminal vertex. The *star of a vertex*  $\text{st}(v)$  is the set of edges emanating from  $v$ , that is

$$\text{st}(v) = \{e : i(e) = v\}. \quad (\text{B.1})$$

In principle, we would like to handle on both the groupoid and geometric dynamics of a given network. One point is that the star of a vertex may be viewed as the combinatorial version of the tangent space to a manifold at a point, rather similar to how the latter may be regarded as an equivalence class of curves through that point. In Bolker et al. (2006) there is defined the notion of a *connection*  $\nabla$  on a graph  $\Gamma$  expressed in terms of a set of one-to-one functions  $\nabla(u, v)$ , one for each oriented edge  $e = (u, v)$  of  $\Gamma$  satisfying the following relationships:

- (1)  $\nabla(u, v) : \text{st}(u) \rightarrow \text{st}(v)$
- (2)  $\nabla(u, v)(u, v) = (v, u)$

$$(3) \nabla(v, u) = (\nabla(u, v))^{-1}$$

Given a graph  $\Gamma$  admits a connection  $\nabla$ , Bolker et al. (2006) define the notion of a *3-geodesic* as a sequence of four vertices  $(u, v, w, z)$  with edges  $\{u, v\}$ ,  $\{v, w\}$  and  $\{w, z\}$  for which

$$\nabla(v, w)(v, u) = (w, z). \quad (\text{B.2})$$

**Remark B.1.** In differential calculus, a ‘connection’ is simply a generalized ‘gradient’ implementing covariant differentiation. We have already encountered a form of this in (5.11). The notion of a graph/network connection introduced here is a more manageable concept, particularly for cognitive modules, and does not involve applying the advanced techniques of calculus.

A *k-geodesic* is defined inductively across a sequence of  $(k + 1)$  vertices. The three consecutive edges  $\{d, e, f\}$  of a 3-geodesic is referred to as *an edge chain*. A *closed geodesic* can then be specified as a sequence of edges  $e_1, \dots, e_n$  such that each consecutive triple  $(e_\alpha, e_{\alpha+1}, e_{\alpha+2})$  is an edge chain for each  $1 \leq \alpha \leq n$ , modulo  $n$ . The geodesic returns to the same pair of edges in the same order. Thus one finds a unique closed geodesic through each pair of edges in the star of the vertex, and as pointed out in Bolker et al. (2006), the set of all closed geodesics completely determines the connection on the graph. Thus for the geometric evolution of our networks, the family  $(\mathbb{G}_A, \nabla_A)$  of local groupoids with connection satisfies:

- (1) Once  $\nabla_A$  is given, then the graph geodesics can be derived iteratively from (B.2).
- (2) Conversely, given the underlying graph of each  $\mathbb{G}_A$ , the connection  $\nabla_A$  is determined by the set of all closed geodesics as specified.

We also have the following useful characterization (Bolker et al., 2006): given  $(\Gamma, \nabla)$ , a subgraph  $\Gamma_0 = (V_0, E_0) \subset \Gamma$  is said to be *totally geodesic* if all geodesics commencing at  $E_0$  remain within  $E_0$ . In other words, for every two adjacent vertices  $u, v$  in  $\Gamma_0$ , we have

$$\nabla(u, v)(\text{st}(u) \cap E_0) \subseteq E_0. \quad (\text{B.3})$$

Note that the above concepts have been formulated graph-theoretically, and as mentioned in Remark B.1, they do not require the usual manipulations of advanced differential calculus.

## B.2 The graph Betti numbers

By analogy with finding the dimensions of the homology groups of a topological manifold, Bolker et al. (2006) specify the notion of *Betti numbers* associated with  $\Gamma$ . This involves the using certain concepts such as an *axial function*  $\varphi$  and *generic direction*  $\xi$ . Thus we regard  $(\Gamma, \nabla)$  as having an axial function  $\varphi$  and write this as  $(\Gamma, \varphi)$  when  $\nabla$  is understood. In which case the *index* of a vertex  $u \in V$  is the number of edges  $e \in \text{st}(u)$  such that the product  $\varphi(e) \cdot \xi < 0$ . Let  $\beta_i(\xi)$  denote the number of vertices  $u$  such that the index at  $u$  is exactly  $i$ . When these values do not depend on the choice of direction  $\xi$ , they are called the *Betti numbers of*  $(\Gamma, \varphi)$ , and satisfy a combinatorial duality condition  $\beta_i(\Gamma, \varphi) = \beta_{d-i}(\Gamma, \varphi)$ , for  $1 \leq i \leq d$ . In certain cases, they can shown to be similar to the indices of a standard Morse function (see Milnor, 1963; Matsumoto, 2001; Bolker et al., 2006) such as  $F_R$  in (3.11). Thus on the underlying graph of the groupoid on which  $F_R$  is defined, we identify  $F_R$  with a Morse function compatible with a generic direction on  $(\Gamma, \varphi)$ . whose index is essentially a measure of the homology of information relay within the graph.; at level  $i$ , we have  $\mu_i = \beta_i(\Gamma, \varphi)$

In fact, to clarify the role of the topological invariants of  $\Gamma$  to those of the surface  $S$ , we need the following description. Firstly,  $S$  taken to be a compact surface permits seeing  $S$  also as a (connected) compact, one-dimensional complex manifold (viz. a Riemann surface) on which a certain analytic group action takes place. The standard way of representing  $\Gamma$  (see e.g. §4 of Bolker et al., 2006) is to identify  $V$  as the (finite) fixed point set, and  $E$  as the (finite) set of one-dimensional orbits of this action. Consequently, the  $\beta_i(\Gamma, \varphi)$  coincide with the usual Betti numbers  $\beta_i(S)$  of  $S$ , and by Poincaré-Hopf we have

$$\chi(S) = \sum_i (-1)^i \beta_i(\Gamma, \varphi). \quad (\text{B.4})$$

## B.3 The graph Laplacian

Suppose now that  $\Gamma = (V, E)$  is an undirected loop-free graph. If the vertices(nodes) are indexed  $1 \leq i \leq N$ , then the *graph Laplacian*  $\Delta$  can be viewed as a symmetric  $N \times N$  matrix defined as follows (see e.g. Biyikoğlu et al., 2007; Tlustý, 2007b):

- (1) If vertices  $(i, j) \in E$  are adjacent, then the corresponding entry in the matrix  $\Delta_{ij} = \Delta_{ji} < 0$ .
- (2) Otherwise,  $\Delta_{ij} = 0$ , and the diagonal terms imply that the sum over rows and columns vanishes, leading to  $\Delta_{ii} = -\sum_{i \neq j} \Delta_{ij}$ .

Note the term ‘weighted Laplacian’ is sometimes used for the operator  $\Delta$ , whereas in other cases ‘Laplacian’ is used for when the negative entries are all  $\Delta_{ij} = -1$ . Specifically, if  $f : V \rightarrow \mathbb{R}$  is a vector function induced by the vertices of  $\Gamma$ , and  $x \sim y$  denotes there is an edge linking  $x$  and  $y$ , then from Biyikoğlu et al. (2007):

$$(\Delta f)(x) = - \sum_{x \sim y} [f(x) - f(y)]. \quad (\text{B.5})$$

Of particular interest are the eigenvalues of  $\Delta$  ordered as  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$  obtainable through the spectrum of an associated operator  $L$ , for which

$$\langle f, Lf \rangle = \sum_{x,y \in V} L_{xy} f(x) f(y) = \sum_{xy \in E} [f(x) - f(y)]^2. \quad (\text{B.6})$$

Also, we have the *Rayleigh Quotients* (Biyikoğlu et al., 2007), given by

$$\begin{cases} \mathcal{R}_\Delta &= \frac{\langle f, \Delta f \rangle}{\langle f, f \rangle} \\ \mathcal{R}_L(f) &= \frac{\sum_{xy \in E} [f(x) - f(y)]^2}{\sum_{x \in V} f(x)^2}. \end{cases} \quad (\text{B.7})$$

In Theorem 5.7 of Dodziuk (1976), estimates on (B.7) leads to showing that, on functions, the eigenvalues of the graph Laplacian converge to those of the continuous Laplacian. Further, in Dodziuk (1976), it is shown that the zeta functions of the former converge to those of the latter, where

$$\zeta^{(n)}(s) = \sum_{\lambda_k^n \neq 0} (\lambda_k^n)^{-s}. \quad (\text{B.8})$$

In the continuous case, sets that are the zero-level sets of the eigenfunctions are called *nodal sets* and *nodal domains* are those sets in which a corresponding eigenfunction takes on one sign and they are separated by nodal sets. Courant’s nodal line theorem (see e.g. Chavel, 1984) states that if the eigenfunctions of a continuous Laplacian on a domain are ordered according to increasing eigenvalues, then the nodes of the  $n$ -th eigenfunction divide the domain into no more than  $n$  nodal domains. In the combinatorial case, for the graph Laplacian, the nodal domains become *sign-graphs*: maximal connected subgraphs on which an eigenfunction carries the same sign. On *weak sign-graphs* the eigenfunction is either  $\geq 0$  or  $\leq 0$ , while on *strong sign-graphs*, the sign of the eigenfunction is either  $> 0$  or  $< 0$ . This leads to an analogue of Courant’s nodal line theorem in the combinatorial case (see e.g. Biyikoğlu et al., 2007): *On a connected graph  $\Gamma$ , the  $n$ -th eigenfunction  $u_n$  of the Laplacian  $\Delta$  admits at most  $n$  weak sign graphs.* The case  $n = 2$  is significant because the corresponding eigenfunction  $u_2$  then splits  $\Gamma$  into exactly two weak sign graphs and  $\lambda_2$  is significant for Brownian motion on the graph and to its first excited energy level.

## C Spin glasses in brief

Spin glass models, as discrete structures, may be based on simplicial decompositions of surfaces usually in some square lattice configuration which can be modified (e.g. from square to triangular). The basic idea leading to the prototypical *2-dimensional Ising model* goes as follows (we follow Surlas (1989) and Carey and Evans (1988)). Firstly, consider a sequence of symbols  $a_i = 0, 1$  and a signal  $v_i$  transmitted across some time interval. Set  $v_i = v$  if  $a_i = 1$ , and  $v_i = -v$  if  $a_i = 0$ . Then let  $a(i, j)$  (for  $1 \leq i \leq m, 1 \leq j \leq m$ ) denote the  $m^2$  bits of information transmitted. These are subject to redundancy relations

$$a(i, m+1) = \sum_{j=1}^m a(i, j) \quad \text{and} \quad a(m+1, j) = \sum_{i=1}^m a(i, j), \quad (\text{C.1})$$

with addition mod 2. The quantity  $m^2/(m+1)^2$  is the rate of the code and measures the redundancy. With noise terms  $y(i, j)$  included, the modified signal is then taken to be  $u(i, j) = v(i, j) + y(i, j)$ . This leads to a simple error-correcting code that is of the Hamming type. Further, the correspondence  $u(i, j) = \frac{1}{2}(\sigma(i, j) + 1)$  between information bits and Ising spins or qubits  $\sigma(i, j)$  in mod 2 addition and spin multiplication respectively, are equivalent.

More specifically, let a qubit  $\sigma(i, j)$  be attached to each edge of some lattice which is to be viewed as a configuration space  $P = \{\pm\}^{\mathbb{Z}^2}$ . On taking  $J_1$  (horizontal) and  $J_2$  (vertical) to be interaction constants, the Hamiltonian  $\mathcal{H}(\sigma)$  is given by

$$\mathcal{H}(\sigma) = - \sum J_1 \sigma(i, j) \sigma(i+1, j) + J_2 \sigma(i, j) \sigma(i, j+1), \quad (\text{C.2})$$

for the appropriate ranges of summation. Suppose we consider  $H(\sigma)$  over a finite lattice given by  $\Lambda_{LM} = \{(i, j) : |i| \leq M, |j| > L\}$ , and then take the thermodynamic limit. If  $J_1, J_2 > 0$ , there are interactions in which the energy is minimized on alignment of all of the spins. Then either

- i) all are  $\uparrow$  or  $\downarrow$ , or,
- ii)  $\sigma(i, j) \equiv 1$ , or  $\sigma(i, j) \equiv -1$ , respectively.

For absolute temperature  $T$ , the equilibrium state is that which minimizes (internal energy)  $-T \cdot$  (entropy). Within the model two competing forces can be realized by the following:

1. One minimizes the internal energy by attempting to align the signs either  $\uparrow$  or  $\downarrow$  to create order: it wins if  $T$  is small.
2. The other, on maximizing entropy, attempts to produce as much chaos as possible: it wins if  $T$  is large.

At finite critical temperature  $T_c$ , chaos wins if  $T \geq T_c$ , and order wins if  $T < T_c$ .

## D References

- Adami, C., Ofria, C., and Collier, T., 2000, Evolution of biological complexity, *Proceedings of the National Academy of Science*, **97**, 4463–4468.
- Adami, C., 2004, Information theory in molecular biology, *Physics of Life Reviews*, **1**, 3-22.
- Ardell, D. H., and Sella, G., 2002, No accident: genetic codes freeze in error-correcting patterns of the standard genetic code, *Phil. Trans. R. Soc. Lond. B* DOI 10.1098/rstb.2002.1071
- Ash, R., 1990, *Information Theory*, Dover Publications, New York.
- Atlan, H., and Cohen, I., 1998, Immune information, self-organization and meaning, *International Immunology*, **10**, 711-717.
- Austin, D. W., Allen, M. S., McCollum, J. M., Dar, R. D., Wilgus, J. R., Saylor, G. S., Samatova, N. F., Cox, C. D., and Simpson, M. L., 2006, Gene network shaping of inherent noise spectra, *Nature*, textbf439, doi:10.1038/nature04194
- Avetisov, V., and Goldanskii, V., 1996, Mirror symmetry breaking at the molecular level, *Proceedings of the National Academy of Science* **93**, 11435–11442.
- Bak, A., Brown, R., Minian, G., and Porter, T., 2006, Global actions, groupoid atlases and related topics, *Journal of Homotopy and Related Structures*, **1**, 1-54.
- Bashford, J. D., Tsohantjis, I., and Jarvis, P. D., 1998, A supersymmetric model for the evolution of the genetic code, *Proceedings of the National Academy of Science*, **95**, 987–992.
- Bashford, J. D., and Jarvis, P. D., 2008, Spectroscopy of the genetic code, in (Abbott, D. et al. eds) *Quantum Aspects of Life*, Imperial College Press, London.
- Belongie, M. L., 1994, Spin glasses and error-correcting codes, *TDA Progress Report 42–118*, 26–36.
- Bennett, C. H., 1982, The thermodynamics of computation: a Review, *International Journal of Theoretical Physics*, **21** (12), 905–940.
- Berger, T., 1971, *Rate Distortion Theory: A mathematical basis for data compression*, Prentice–Hall Inc., Englewood Cliffs, NJ.
- Bertman, M. O., and Jungck, J. R., 1978, Some unresolved mathematical problems in genetic coding, *Notices of the Amer. Math. Soc.*, **25** A-174.
- Biyikoğlu, T., Leydold, J., and Stadler, P. F., 2007, Laplacian Eigenvectors of Graphs, *Lecture Notes in Math.* **1915**, Springer-Verlag, Berlin Heidelberg.
- Bolker, E. D., Guillemin, V. W., & Holm, T. S., 2006, How is a graph like a manifold?, to appear. <http://arxiv.math.CO/0206103>
- Bos, R., 2010, Continuous representations of groupoids, *Houston J. Math.*, to appear. [arXiv:math/0612639](http://arXiv:math/0612639)
- Brown, R., 2006, *Topology and Groupoids*, BookSurge LLC.
- Burago, D., Burago, Y., and Ivanov, S., 2001, *A Course in Metric Geometry*, American Mathematical Society, Providence, RI.

- Carey, A., and Evans, D. E., 1988, The operator algebras of the two-dimensional Ising model, in *'Braids' (Santa Cruz, CA, 1986)*, 117–165, *Contemp. Math.* **78**, Amer. Math. Soc., Providence, RI.
- Chavel, I., 1984, *Eigenvalues in Riemannian Geometry*, Academic Press, New York.
- Ciliberti, S. Martin, O., and Wagner, A., 2007, Innovation and robustness in complex regulatory genetic networks, *Proceedings of the National Academy of Science* **104**, 13591–13596.
- Cohen, I., 2000, *Tending Adam's Garden: Evolving the Cognitive Immune Self*, Academic Press, New York.
- Cohen, I., and Harel D., 2007, Explaining a complex living system: dynamics, multiscaling and emergence, *Journal of The Royal Society Interface* **4**, 175–182.
- Connes, A., 1994, *Noncommutative Geometry*, Academic Press, San Diego, CA.
- Cover, T., and Thomas, J., 1991, *Elements of Information Theory*, John Wiley and Sons, New York.
- Crick, F., 1966, Codon-anticodon pairing: the wobble hypothesis, *J. Mol. Biol.* **19**, 548–553.
- Crick, F., 1968, The origin of the genetic code, *J. Mol. Biol.* **38**, 367–379.
- Crooks, G. E. and Brenner, S., 2004, Protein secondary structure: entropy, correlations and prediction, *Bioinformatics* **20**(10), 1603–1611.
- Dawkins, R., 1976, *The Selfish Gene*, Oxford Univ. Press, London.
- Dodziuk, J., 1976, Finite-difference approach to the Hodge theory of harmonic forms, *Amer. J. Math* **98**(1), 79–104.
- Dretske, F., 1981, *Knowledge and the Flow of Information*, MIT Press, Cambridge, MA.
- Feynman, R., 1996, *Feynman Lectures on Computation*, Addison-Wesley, Reading, MA.
- Franzosi, R., and M. Pettini, 2004, Theorem on the origin of phase transitions, *Physical Review Letters*, **92**:060601.
- Gent, I. P., Kelsey, T., Linton, S., Pearson, J., and Roney-Dougal, C. M., 2010, Groupoids and conditional symmetry, preprint, University of St. Andrews, UK.
- Glazebrook, J. F., and Wallace, R., 2009a, Small worlds and red queens in the global workspace: an information-theoretic approach, *Cognitive Systems Research*, **10**, 333–365.
- Glazebrook, J. F., and Wallace, R., 2009b, Rate distortion manifolds as model spaces for cognitive information, *Informatica* **33** (2009), 309–345.
- Glazebrook, J. F., and Wallace, R., 2010, Rate distortion coevolutionary dynamics and the flow nature of cognitive epigenetic systems, submitted.
- Golubitsky, M., and Stewart, I., 2006, Nonlinear dynamics and networks: the groupoid formalism, *Bulletin of the American Mathematical Society*, **43**, 305–364.
- Gupta, M. K., 2006, The quest for error correction in biology, *IEEE Engineering in Medicine and Biology Magazine*, 46–53.
- Hornos, J. E., and Hornos, Y. M., 1994, A search for symmetries in the genetic code, *Journal of Biological Physics*, **20**, 289–294.
- Jiménez-Montaña, M. A., de la Mora-Basáñez, C. R., and Pöschel, T., 1996, The hypercube structure of the genetic code explains non-conservative aminoacid substitutions *in vivo* and *in vitro*, *BioSystems*, **39**, 117–125.
- Jukes, T. H., 1983, Evolution of the amino acid code, pp. 191–207 in (Nei, M. et al. eds) *Evolution of Genes and Protein*, Sinauer, Sunderland, MA.
- Khinchin A., 1957, *The Mathematical Foundations of Information Theory*, Dover Publications, New York.
- Koonin, E., and Novozhilov, A., 2009, Origin and evolution of the genetic code: the universal enigma, *Life*, **61**, 99–111.
- Kurzynski, M., 2006, *The Thermodynamic Machinery of Life*, Springer-Verlag, Berlin Heidelberg New York.
- Landau, L., and Lifshitz E., 2007, *Statistical Physics (I)* (3rd Ed.), Elsevier, New York.
- Lee, J., 2000, *Introduction to Topological Manifolds*, Springer, New York.
- Levinthal, L., 1968, Are there pathways for protein folding?, *J. Chim. Phys. PCB*, **65**, 44–45.

- Lisewski, A. M., 2008, Random amino acid mutations and protein misfolding lead to Shannon limit in sequence-structure communication, *PloS ONE*, **3**(9), e3110 doi:10.371/journal.pone.0003110
- Matsumoto, Y., 2001, *An introduction to Morse Theory*, Translations of the American Mathematical Society **208**, Providence, RI.
- McEliece, R. J., 2004, *The Theory of Information and Coding*, Encyclopedia of Mathematics and its Applications, Vol. 86, Cambridge University Press.
- Milnor, J., 1963, *Morse Theory*, Princeton University Press, Princeton, NJ.
- Pettini, M., 2007, *Geometry and Topology in Hamiltonian Dynamics*, Springer, New York.
- Protter, P., 1995, *Stochastic Integration and Differential Equations: A New Approach*, Springer, New York.
- Prügel-Bennett, A. and Shapiro, J. L., 1994, Analysis of genetic algorithms using statistical mechanics, *Physical Review Letters*, **72** no. 9., 1305–1309.
- Ringel, G., and J. Youngs, 1968, Solutions of the Heawood map-coloring problem, *Proceedings of the National Academy of Sciences*, **60**, 438-445.
- Rose, K., Gurewitz, E., and Fox, G. C., 1990, Statistical mechanics and phase transitions in clustering, *Physical Review Letters*, **65** No. 6, 945–948.
- Schrödinger, E., 1967, *What is Life?*, Cambridge University Press.
- Sella, G. and Ardell, D. H., 2002, The impact of message mutation on the fitness of the genetic code, *J. Mol. Evol.* **54**, 638–651.
- Sella, G. and Ardell, D. H., 2006, The coevolution of genes and genetic codes: Crick’s frozen accident revisited, *J. Mol. Evol.* **63**, 297–313.
- Söll, D., and RajBhandary, U. L., 2006, The genetic code—Thawing the ‘frozen accident’, *J. Biosci.* **31**(4), 459–463.
- Sourlas, N., 1989, Spin-glass models as error-correcting codes *Nature* **339**, 693 - 695  
doi:10.1038/339693a0
- Stewart, I., Golubitsky, M., and Pivato, M., 2003, Symmetry groupoids and patterns of synchrony in coupled cell networks, *SIAM Journal of Applied Dynamical Systems*, **2**, 609-646.
- Stewart, I., 1994, Broken symmetry in the genetic code?, *New Scientist* **1915**, 16.
- Trusty, T., 2007a, A model for the emergence of the genetic code as a transition in a noisy channel, *Journal of Theoretical Biology* **249**, 331–342.
- Trusty, T., 2007b, A relation between the multiplicity of the second eigenvalue of a graph Laplacian, Courant’s nodal line theorem and the substantial dimension of tight polyhedral surfaces, *Electric Journal of Linear Algebra* **16**, 315–324.
- Trusty, T., 2008a, Rate-distortion scenario for the emergence and evolution of noisy molecular codes, *Physical Review Letters* **100**, 048101 (1-4).
- Trusty, T., 2008b, A simple model for the evolution of molecular codes driven by the interplay of accuracy, diversity and cost, *Physical Biology*, **5**, 016001.
- Vetsigian, K., Woese, C., and Goldenfeld, N., 2006, Collective evolution and the genetic code, *Proceedings of the National Academy of Science* **103**, 10696-10701.
- Wallace, R., 2005, *Consciousness: A Mathematical Treatment of the Global Neuronal Workspace Model*, Springer, New York.
- Wallace, R., 2010a, A rate distortion approach to protein symmetry, *BioSystems*, to appear.
- Wallace, R., 2010b, Metabolic constraints on the evolution of genetic codes: Did multiple ‘preaerobic’ ecosystem transitions entrain richer dialects via Serial Endosymbiosis  
<http://precedings.nature.com/documents/4120/version/3>.
- Wallace, R., 2010c, Many *in vivo* ‘protein folding codes’ can be inferred from empirical classifications using Trusty’s topological approach, to appear.
- Wallace, R., and Wallace R.G., 1998, Information theory, scaling laws, and the thermodynamics of evolution, *Journal of Theoretical Biology*, **192**, 545-559.

Wallace, R., and Wallace R.G., 1999, Organisms, organizations and interactions: an information theory approach to biocultural evolution, *BioSystems*, **51**, 101-119.

Wallace, R., and Wallace, D., 2009, *Gene Expression and its Discontents: The social production of pandemic chronic disease*, Springer, New York.

Wallace, R., and Wallace, D., 2008, Punctuated equilibrium in statistical models of generalized coevolutionary resilience: how sudden ecosystem transitions can entrain both phenotype expression and Darwinian selection, *Transactions on Computational Systems Biology IX*, LNBI 5121: 23 –85.

Wallace, R., 2009, Cultural epigenetics: on the heritability of complex diseases, *Transactions on Computational Systems Biology*, to appear.

Wallace, R., and Fullilove M., 2008, *Collective Consciousness and Its Discontents: Institutional distributed cognition, racial policy, and public health in the United States*, Springer, New York.

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

Yockey, H. P., 2005, *Information Theory, Evolution and the Origin of Life*, Cambridge University Press.